

1

Version dated: July 4, 2023

2 RH: Deep Learning and Phylogeography

3 [EIC] Deep learning and likelihood approaches for viral
4 phylogeography converge on the same answers whether
5 the inference model is right or wrong

6 AMMON THOMPSON^{1,*}, BENJAMIN LIEBESKIND², ERIK J. SCULLY², MICHAEL LANDIS^{3,*}

7 ¹*Participant in an education program sponsored by U.S. Department of Defense (DOD) at the
National Geospatial-Intelligence Agency, Springfield, VA, 22150, USA*

8 ²*National Geospatial-Intelligence Agency, Springfield, VA, 22150, USA*

9 ³*Department of Biology, Washington University in St. Louis, Rebstock Hall, St. Louis, Missouri,
11 63130, USA*

12 *Corresponding authors: E-mail: Ammon.M.Thompson.ctr@nga.mil and

13 michael.landis@wustl.edu.

¹⁴ *Abstract.*— Analysis of phylogenetic trees has become an essential tool in epidemiology.

¹⁵ Likelihood-based methods fit models to phylogenies to draw inferences about the

¹⁶ phylodynamics and history of viral transmission. However, these methods are

¹⁷ computationally expensive, which limits the complexity and realism of phylodynamic

¹⁸ models and makes them ill-suited for informing policy decisions in real-time during rapidly

¹⁹ developing outbreaks. Likelihood-free methods using deep learning are pushing the

²⁰ boundaries of inference beyond these constraints. In this paper, we extend, compare and

²¹ contrast a recently developed deep learning method for likelihood-free inference from trees.

²² We trained multiple deep neural networks using phylogenies from simulated outbreaks that

²³ spread among five locations and found they achieve [R2.1] close to the same levels of

²⁴ accuracy as Bayesian inference under the true simulation model. We compared robustness

²⁵ to model misspecification of a trained neural network to that of a Bayesian method. We

²⁶ found that both models had comparable performance, converging on similar biases. [EIC,

²⁷ R1.2] We also implemented a method of uncertainty quantification called conformalized

²⁸ quantile regression which we demonstrate has similar patterns of sensitivity to model

²⁹ misspecification as Bayesian highest posterior intervals (HPI) and greatly overlap with

³⁰ HPIs, but have lower precision (more conservative). Finally, we trained and tested a neural

³¹ network against phylogeographic data from a recent study of the SARS-CoV-2 pandemic in

³² Europe and obtained similar estimates of region-specific epidemiological parameters and

³³ the location of the common ancestor in Europe. Along with being as accurate and robust

³⁴ as likelihood-based methods, our trained neural networks are on average over 3 orders of

³⁵ magnitude faster. Our results support the notion that neural networks can be trained with

³⁶ simulated data to accurately mimic the good and bad statistical properties of the

³⁷ likelihood functions of generative phylogenetic models.

³⁸ (Keywords: phylogeography, SSE, phylodynamics, machine learning, deep learning,

³⁹ epidemiology)

INTRODUCTION

41 Viral phylodynamic models use genomes sampled from infected individuals to [Authors]
42 infer the evolutionary history of a pathogen and its spread through a population (Holmes
43 and Garnett 1994; Volz et al. 2013). By linking genetic information to epidemiological
44 data, such as the location and time of sampling, these generative models can provide
45 valuable insights into the transmission dynamics of infectious diseases, especially in the
46 early stages of cryptic disease spread when it is more difficult to detect and track (Holmes
47 et al. 1995; Rambaut et al. 2008; Lemey et al. 2009; Pybus et al. 2012; Worobey et al.
48 2016, 2020; Lemey et al. 2021; Washington et al. 2021; Pekar et al. 2022). This information
49 can be used to inform public health interventions and improve our understanding of the
50 evolution and spread of pathogens. Many phylodynamic models are adapted from
51 state-dependent birth-death (SDBD) processes or, equivalently, state-dependent
52 speciation-extinction (SSE) models (Maddison et al. 2007; FitzJohn 2012; Kühnert et al.
53 2014; Beaulieu and O’Meara 2016). [AE] These birth-death models correspond to the
54 well-known Susceptible-Infectious-Recovered (SIR) model during an exponential growth
55 phase, when nearly all individuals in the population are susceptible to infection (Anderson
56 and May 1979). The simplest SIR models only track the number of susceptible, infected,
57 and recovered individuals across populations over time, with more advanced models also
58 allowing the movement of individuals among localized populations. The phylodynamic
59 models we are interested in track the incomplete transmission tree (phylogeny) of sampled,
60 infected individuals that emerges from host-to-host pathogen spread among populations
61 over space and time. Within this broader context, we will refer to the state as location and
62 the models as location-dependent birth-death (LDBDS) models that include serial
63 sampling of taxa (Kühnert et al. 2016).

64 Analysts [Authors] typically fit these birth-death models to data using
65 likelihood-based inference methods, such as maximum likelihood (Maddison et al. 2007;
66 Richter et al. 2020) or Bayesian [Authors] inference (Kühnert et al. 2016; Scire et al. 2020).

67 Likelihood-based inference relies upon a likelihood function to evaluate the relative
68 probability (likelihood) that a given phylogenetic pattern (i.e., topology, branch lengths,
69 and tip locations) was generated by a phylodynamic process with particular model
70 parameter values. In this sense the likelihood of any possible phylodynamic data set is
71 mathematically encoded into the likelihood as a function of (unknown) data-generating
72 model parameters.

73 Computing the likelihood requires high-dimensional integration over a large and
74 complex space of evolutionary histories. Analytically integrated likelihood functions,
75 however, are not known for LDBDS models. Methods developers instead use ordinary
76 differential equation (ODE) solvers (Maddison et al. 2007; Kühnert et al. 2016) to
77 numerically approximate the integrated likelihood. These clever approximations perform
78 well statistically, but are too computationally expensive to use with large epidemic-scale
79 data sets. Thus, while Nextstrain (Hadfield et al. 2018) and similar efforts have provided
80 useful visualizations to policy makers during the COVID response, most phylogeographical
81 methods are used forensically, providing insight on the past, and are not used to provide
82 parameter estimates in response to emerging events to inform policy decisions in real-time
83 due to the complexity and long run-times of these models.

84 As phylodynamic models become more biologically realistic, they will necessarily
85 grow more mathematically complex, and therefore less able to yield likelihood functions
86 that can be approximated using ODE methods. Because of this, phylodynamic model
87 developers tend to explore only models for which a likelihood-based inference strategy is
88 readily available. As a consequence, [Authors] the lack of scalable inference methods
89 impedes the design, study, and application of richer phylodynamic models of disease
90 transmission[EIC], in particular, and richer phylogenetic models of lineage diversification,
91 in general.

92 To avoid the computational limitations associated with likelihood-based methods,
93 deep learning inference methods that are likelihood-free have emerged as a complementary

94 framework for fitting a wide variety of evolutionary models (Bokma 2006). Deep learning
95 methods rely on training many-layered neural networks to extract information from data
96 patterns. These neural networks can be trained with simulated data as another way to
97 approximate the latent likelihood function (Cranmer et al. 2020). Once trained, neural
98 networks have the benefit of being fast, easy to use, and scalable. Recently, likelihood-free
99 deep learning neural network methods have successfully been applied to phylogenetics
100 (da Fonseca et al. 2020; Suvorov et al. 2020; Nesterenko et al. 2022; Solis-Lemus et al.
101 2022; Suvorov and Schrider 2022) and phylodynamic inference (Lambert et al. 2022;
102 Voznica et al. 2022).

103 Here we extend new methods of deep learning from phylogenetic trees (Lambert
104 et al. 2022; Voznica et al. 2022) to explore their potential when applied to phylogeographic
105 problems in geospatial epidemiology. Phylodynamics of birth-death-sampling processes
106 that include migration among locations have been under development for more than a
107 decade (Stadler 2010; Stadler et al. 2012; Kühnert et al. 2014, 2016; Scire et al. 2020; Gao
108 et al. 2022, 2023). Given the added complexity of location-specific dynamics (e.g.
109 location-specific infection rates) and recent successes in deep learning with phylogenetic
110 time trees (Voznica et al. 2022) under state-dependent diversification models (Lambert
111 et al. 2022), we sought to evaluate this approach when applied to viral phylodynamics and
112 phyogeography by including location data when training deep neural networks with
113 phylogenetic trees.

114 A current limitation of likelihood-free approaches is that it remains unknown how
115 brittle the inference machinery is when the assumptions used for simulation and training
116 are violated (Schmitt et al. 2022). For example, a brittle deep learning method would be
117 more easily misled by model misspecification when compared to a likelihood-based method.
118 Likelihood approaches may have some advantages because the simplifying assumptions are
119 explicit in the likelihood function while for trained neural networks it is difficult to know
120 how those [R2.6] same assumptions implemented in the simulation are encoded in data

121 patterns in the training data and learned network weights. However, with complex
122 likelihood models, there may be unexpected interactions among simplifying assumptions
123 that can result in large biases when applied to real-world data (Gao et al. 2023).
124 Characterizing the relative robustness and brittleness of these two inference paradigms is
125 essential for those who wish to confidently develop and deploy likelihood-free methods of
126 inference from real world data.

127 To explore relative robustness to model misspecification, we trained multiple deep
128 convolutional neural networks (CNNs) with transmission trees generated from epidemic
129 simulations. [R2.1] We were able to achieve accuracy very close to that of a
130 likelihood-based approach and through several model misspecification experiments show
131 that our CNNs are no more sensitive to model violations than the likelihood approach.
132 Significantly, both methods consistently show similar biases induced by model violations in
133 test data sets. We find that for the models tested here, the migration rate estimates are
134 highly sensitive to misspecification of infection rate and sampling rates, but that estimates
135 of the infection and sampling rates are fairly robust to misspecification of the migration
136 models. We also show that the rate parameter estimates are fairly robust to
137 misspecification of both the number of locations in the model and phylogenetic error. [EIC,
138 R1.2] We also estimated prediction intervals for the rate parameters and compared and
139 contrasted their performance to the Bayesian highest posterior density intervals (HPI). We
140 show that they produce intervals that greatly overlap with HPIs in all experiments, but
141 have, on average, wider intervals making them relatively conservative. Finally, we
142 compared a simulation-trained neural network to a recent phylodynamic study of the first
143 wave of the COVID pandemic in Europe (Nadeau et al. 2021) and obtain similar inferences
144 about the dynamics and history of SARS-CoV-2 in the European clade.

145

METHODS

146 First, we define the SIR model we assume here that is approximately equivalent to

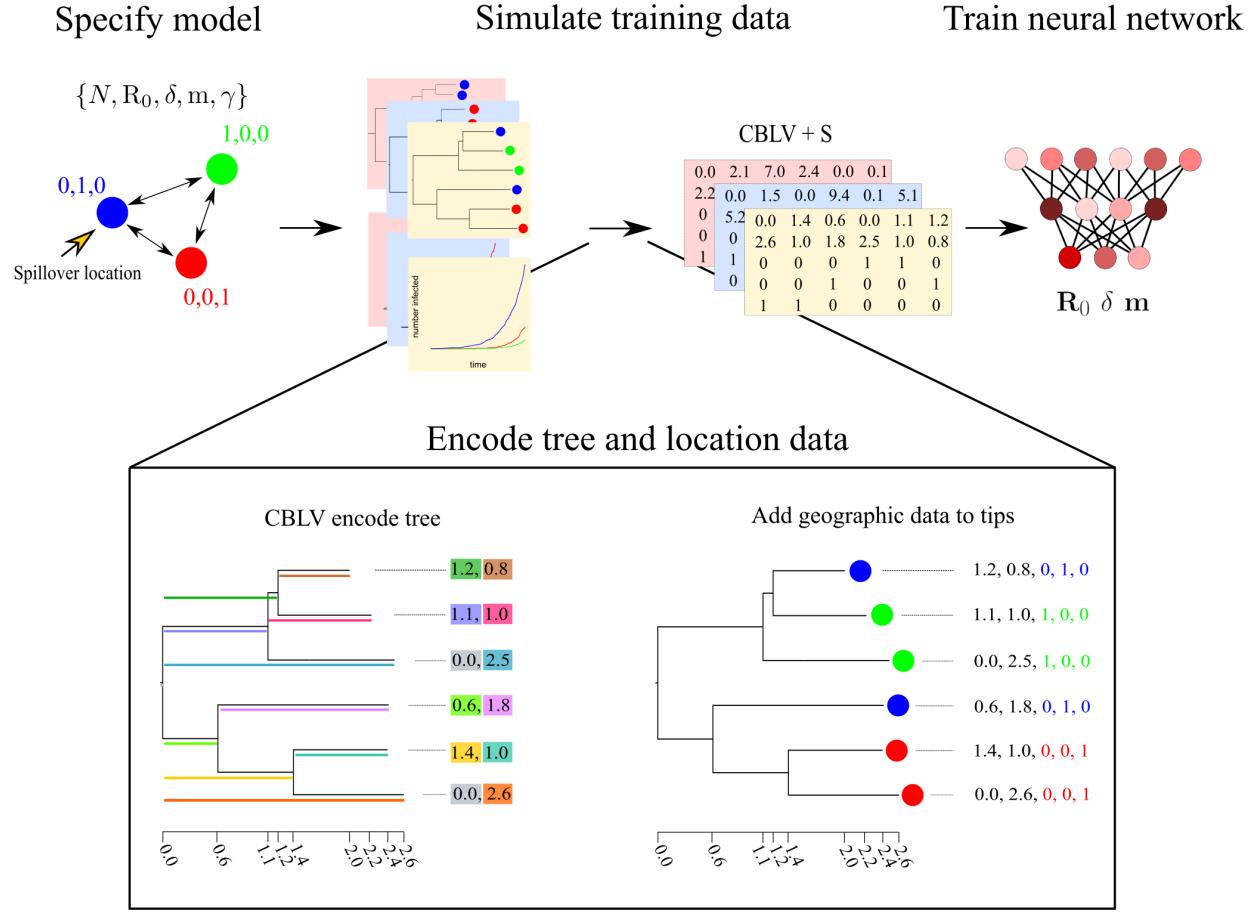


Figure 1: [AE.5] Simulation and tree encoding pipeline for generating training data. 1) Specify a model, for example an SIR model with serial sampling and migration among three locations (colored circles). 2) Run simulations of outbreaks under the model to generate population trajectories and phylogenetic trees. 3) encode trees and location data into the Compact Bijective Ladderized Vector + States (CBLV+S) format. 4) Train the neural network with CBLV+S training data.

147 the LDBDS model (Kühnert et al. 2016). Following that, is a description of the simulation
 148 method to generate the training, validation, and test data sets of phylogenies under the
 149 model. [AE] The simulation and data processing pipeline is shown in Figure 1. We next
 150 describe our implementation of simulation-trained deep learning inference with
 151 convolutional neural networks (CNN) as well as a likelihood-based method using Bayesian
 152 inference. We then describe our methods for measuring and comparing their performance
 153 when tested against data sets generated by simulations under the inference model as well
 154 as several data sets simulated under models that violate assumptions of the inference
 155 model. Finally, we describe how we tested our simulation-trained CNN against a real-world
 156 data set.

157 *Model definition*

158 We first define a general location-dependent SIR stochastic process used for simulations
 159 and likelihood function derivation in the format of reaction equations we specified in
 160 MASTER (Vaughan and Drummond 2013). Reaction equations 1 through 4 specify the
 161 SIR compartment model with migration and serial sampling where S , I , and R denote the
 162 number of individuals in each compartment. The S and I compartments are indexed by
 163 geographic location using i and j . N_i is the total population size in location i and
 164 $N_i = S_i + I_i + R_i$. [Authors] To simplify notation, we consider all local recoveries to lead to
 165 the same global compartment and absorbing state, R . The symbols for each rate parameter
 166 is placed above each reaction arrow.



₁₆₇ We parameterize the model with the basic reproduction number in location i , R_{0_i} ,
₁₆₈ which is related to β_i and δ_i by equation 5,

$$R_{0_i} = \frac{\beta_i}{\gamma + \delta_i}. \quad (5)$$

₁₆₉ In particular, our study considers a location-independent SIR (LISIR) model with
₁₇₀ sampling that assumes R_{0_i} was equal among all locations, and a location-dependent
₁₇₁ (LDSIR) model with sampling that assumes R_{0_i} varied among locations. During the
₁₇₂ exponential growth phase of an outbreak, the LISIR and LDSIR models are equivalent to
₁₇₃ the location-independent birth-death-sampling (LIBDS) and location-dependent
₁₇₄ birth-death-sampling (LDBDS) models, respectively, that are often used in viral
₁₇₅ phylogeography (Kühnert et al. 2014, 2016; Douglas et al. 2021).

₁₇₆ Each infectious individual transitions to recovered at rate γ . We assumed that
₁₇₇ sampling a virus in an individual occurs at rate δ_i in location i and immediately removes
₁₇₈ that individual from the infectious compartment and places them in the recovered
₁₇₉ compartment. Thus the effective recovery rate in location i is $\gamma + \delta_i$. The above reactions
₁₈₀ correspond to the following coupled ordinary differential equations.

$$\begin{aligned} \frac{dS_i}{dt} &= -\frac{\beta_i}{N_i} S_i I_i \\ \frac{dI_i}{dt} &= \frac{\beta_i}{N_i} S_i I_i + \sum_{j \neq i}^n m_{ij} I_j - \sum_{j \neq i}^n m_{ji} I_i - (\gamma + \delta_i) I_i \\ \frac{dR}{dt} &= \sum_{i=1}^n (\gamma + \delta_i) I_i \end{aligned} \quad (6)$$

₁₈₁ When the migration rate is constant among locations and the model is a
₁₈₂ location-independent SIR model, or equivalently, LIBDS, [authors:] and we set
₁₈₃ $S_i(t = 0) \approx N_i$ at the beginning of the outbreak, the equation set 6 reduces to

[authors:]

$$\begin{aligned}\frac{dS_i}{dt} &= -\beta I_i \\ \frac{dI_i}{dt} &= \beta I_i + m \left(\sum_{j \neq i}^n I_j - (n-1)I_i \right) - (\gamma + \delta)I_i \\ \frac{dR}{dt} &= (\gamma + \delta) \sum_{i=1}^n I_i\end{aligned}$$

184

185 The number of infections and the migration of susceptible individuals is at
186 negligible levels on the timescales investigated here. The infection rate is, therefore,
187 approximately constant and the migration of susceptible individuals can be safely ignored
188 requiring only migration of infectious individuals to be simulated.

189 At the beginning of an outbreak, it is often easier to know the recovery period from
190 clinical data than the sampling rate which requires knowing the prevalence of the disease.
191 Therefore, we treat the average recovery period as a known quantity and use it to make the
192 other two parameters (the sampling rate and the basic reproduction number R_0)
193 identifiable. This was done by fixing the corresponding rate parameter in the likelihood
194 function to the true simulated value for each tree, and by adding the true simulated value
195 to the training data for training the neural network.

196

Simulated training and validation data sets

197 Epidemic simulations of the SIR+migration model that approximates the LIBDS process
198 were performed using the MASTER package v. 6.1.2 (Vaughan et al. 2014) in BEAST 2 v.
199 2.6.6 (Bouckaert et al. 2019). MASTER allows users to simulate phylodynamic data sets
200 under user-specified epidemiological scenarios, for which MASTER simultaneously
201 simulates the evolution of compartment (population type) sizes and tracks the branching
202 lineages (transmission trees in the case of viruses) from which it samples over time. We

203 trained neural networks with these simulated data to learn about latent populations from
 204 the shape of sampled and subsampled phylogenies. In addition to the serial sampling
 205 process, at the end of the simulation 1% of infected lineages were sampled. In MASTER
 206 this was approximated by setting a very high sampling rate and very short sampling time
 207 such that the expected number sampled was approximately 1%. This final sampling event
 208 was required to make a 1-to-1 comparison of the likelihood function used for this study (see
 209 Likelihood method description below) which assumes at least one extant individual was
 210 sampled to end the process. Coverage statistics from our MCMC samples closely match
 211 expectations (see Likelihood method description below; SI Figure 2 C). Simulation
 212 parameters under LIBDS and LDBDS models for training the neural network under the
 213 phyogeography model were drawn from the following distributions:

$$\begin{aligned}
 R_0 &\sim \text{Uniform}(2, 8) \\
 \delta &\sim \text{Unif}(0.0001, 0.005) \\
 m &\sim \text{Uniform}(0.0001, 0.005) \\
 \gamma &\sim \text{Unif}(0.01, 0.05)
 \end{aligned} \tag{7}$$

[authors:] spillover location $\sim \text{Multinomial}(k = 1, p_i = 1/5)$, for 5 locations

214 All five locations had initial population sizes of 1,000,000 susceptible individuals and
 215 one infected individual in a randomly sampled spillover location. Simulations were run for
 216 100 time units or until 50,000 individuals had been infected to restrict simulations to the
 217 approximate exponential phase of the outbreak. For the experiments comparing the CNN
 218 to the likelihood-based method under the LIBDS model, if this population threshold was
 219 reached the simulation was rejected. This criterion was not enforced for simulations under
 220 the LDBDS model. This ensured the LIBDS model used in the likelihood-based analyses
 221 are equivalent to more complex density-dependent SIR models. After simulation, trees with

222 500 or more tips were uniformly and randomly downsampled to 499 tips and the sampling
223 proportion was recorded for training the neural networks and to adjust estimates of δ .

224 We simulated 410,000 outbreaks under these LIBDS settings to generate the
225 training, validation, and test sets for deep learning. Any simulation that generated a tree
226 with less than 20 tips was discarded, leaving a total of 111,157 simulated epidemiological
227 data sets. Of these, 104,157 data sets were used to train and 7,000 were used to validate
228 and test each CNN. A total of 193,110 LDBDS data sets were simulated, with 186,110 used
229 to train and 7,000 used to validate and test the LDBDS CNNs.

230 [AE.3] To make phylodynamic inferences about the first wave of the SARS-CoV-2
231 epidemic in Europe we used the LDBDS model on the data set from Nadeau et al. (2021).
232 Training simulation parameters for the LDBDS process were drawn from the same
233 distributions as LIBDS except R_0 which was unique for each location. We assume that the
234 variability of R_0 among different pathogens (simulated outbreaks) is greater than the
235 variability of the same pathogen's R_0 among different locations within the same simulation.
236 To implement this assumption, all R_0 was drawn from a joint distribution to narrow the
237 magnitude of differences among locations within simulations to be within 6 of each other
238 but expand the magnitude of differences between simulations to range from 0.9 to 15:

$$\alpha \sim \text{Uniform}(3.9, 12)$$

$$R_{0_i} | \alpha \sim \text{Uniform}(\alpha - 3, \alpha + 3)$$

239 For the empirical analysis, population sizes at each location were also set to 500,000
240 and instead of running the simulations for 100 time units, time was scaled by the recovery
241 period, $1/\gamma$, and was drawn from a uniform distribution:

time $\sim \text{Uniform}(1, 20)$

242 *Simulated test data sets with and without model misspecification*

243 [authors] All simulation models used for training and testing are listed in Table 1.

244 We first simulated a test set of 138 trees under the training model to compare the accuracy
245 of the CNN and the likelihood-based estimates when the true model is specified. These
246 data sets were simulated by random draws of parameter values from the same distributions
247 described above for generating the training data set.

248 Sensitivity to model misspecification for each of the three rate parameters, R_0 , δ ,
249 and m , was tested. All sensitivity experiments used the same LIBDS model for inference
250 for both the CNN and the Likelihood-based methods. Sensitivity experiments were
251 conducted by simulating a test data set of trees that were generated by an epidemic
252 process that was more complex than or different from the LIBDS model.

253 The tree data set for the misspecified R_0 experiment consisted of simulating
254 outbreaks where each location had a unique R_0 drawn from the same distribution as above.
255 Likewise, the misspecified sampling model test set was generated by simulating outbreaks
256 where each location had a unique sampling rate, δ , drawn from the same distribution used
257 for the global sampling rate described above. For the misspecified migration model, a
258 random pair of coordinates, each drawn from a uniform(0,5) distribution in a plane, were
259 generated for the five locations, and a pairwise migration rate was computed such that
260 pairwise migration rates were symmetric and proportional to the inverse of their euclidean
261 distances and the average pairwise migration rate was equal to a random scalar which was
262 also drawn from a uniform distribution (see equations 7 above).

263 The tree set for the misspecified number of locations experiment was generated by

Description	Simulation model parameters and data
Generate training data	$\{N, R_0, \delta, m, \gamma, \Psi\}$
Misspecify R_0	$\{N, R_{01}, R_{02}, R_{03}, R_{04}, R_{05}, \delta, m, \gamma, \Psi\}$
Misspecify δ	$\{N, R_0, \delta_1, \delta_2, \delta_3, \delta_4, \delta_5, m, \gamma, \Psi\}$
Misspecify m	$\{N, R_0, \delta, m_{ij} \forall i \neq j \in \{1, \dots, N\}, \gamma, \Psi\}$
Misspecify number of locations	$\{2N, R_0, \delta, m, \gamma, \Psi\}$
Tree error	$\{N, R_0, \delta, m, \gamma, \Psi^{\text{error}}\}$
Analyze Nadeau et al. (2021) dataset	$\{N, R_{01}, R_{02}, R_{03}, R_{04}, R_{05}, \delta, m, \gamma, \Psi\}$

Table 1: [authors] Models used in this study. All simulations assume an SIR compartmental epidemic model. $N = 5$ is the number of locations, R_0 is the basic reproduction number, δ is the sampling rate, m is the migration rate, γ is the recovery rate (treated as data), and Ψ is the phylogenetic tree + locations (also treated as data).

264 simulating outbreaks among ten locations instead of five. After simulations, six locations
 265 were chosen at random and re-coded as being sampled from the same location.

266 To generate a test set where the time tree used for inference has incorrect topology
 267 and branch lengths, we implemented a basic pipeline of tree inference from simulated
 268 genetic data to mimic a worst case real world scenario. We simulated trees under the same
 269 settings as before. Phylogenetic error was introduced in two ways: the amount of site data
 270 (short sequences) and misspecification of the DNA sequence evolution inference model
 271 [Authors:] using seq-gen V. 1.3.2 (Rambaut and Grassly 1997). We simulated the evolution
 272 of a 200 base-pair sequence under an HKY model with $\kappa = 2$, equal base frequencies and 4
 273 discretized-gamma(2, 2) rate categories for among site rate variation. The simulated
 274 alignment as well as the true tip dates (sampling times) was then used to infer test trees.
 275 Test tree inference was done using IQ-Tree v. 2.0.6 (Minh et al. 2020) assuming a
 276 Jukes-Cantor model of evolution where all transition rates are equal. The inference model
 277 also assumed no among-site rate variation. The number of shared branches between the
 278 true transmission tree and the test tree inferred by IQ-Tree was measured using gotree v.
 279 0.4.2 (Lemoine and Gascuel 2021). Polytomies were resolved using phytools (Revell 2012)
 280 and a small, random number was added to each resolved branch. These trees were then
 281 used for likelihood inference and CNN prediction.

282

Deep learning inference method

283 The resulting trees and location metadata generated by our pipeline were converted to a
284 modified CBLV format (Voznica et al. 2022), which we refer to as the CBLV+S (+State of
285 character, *e.g.* location) format (Figure 1). The CBLV format uses an in-order tree
286 traversal to translate the topology and branch lengths of the tree into an $2 \times n$ matrix
287 where n is the maximum number of tips allowed for trees. [authors] The matrix is
288 initialized with zeroes. We then fill the matrix starting with the root then proceed to the
289 tip with largest root-to-tip distance rather than starting with that tip as in Voznica et al.
290 (2021). We chose this to separate the the zero value of the root age from the zeroes used to
291 pad matrices where the tree has less than the maximum number of tips, though we expect
292 this to make marginal to no difference in performance. The CBLV representation gives
293 each sampled tip a pair of coordinates in ‘tree-traversal space’. Our CBLV+S format
294 associates geographic information corresponding with each sampled taxon by appending
295 each vector column with a one-hot encoding vector of length g states to yield a $(2 + g) \times n$
296 CBLV+S matrix. The CBLV+S format allows for multiple characters and/or states to be
297 encoded, extending the single binary character encoding format introduced by Lambert
298 et al. (2022). Our study uses CBLV+S to encode a single character with $g = 5$
299 location-states. In addition to the the CBLV+S data, we also include a few tree summary
300 statistics and known simulating parameters; [authors] the number of tips, mean branch
301 length, the tree height and the recovery rate and the subsampling proportion. Trees were
302 rescaled such that their mean branch length was the default for phylodeep (Voznica et al.
303 2022) before training and testing of the CNN. The mean pre-scaling branch length and tree
304 heights were also fed into the neural networks. Trees were not rescaled for the
305 likelihood-based analysis. Recall that tree height did not vary for the LIBDS CNN training
306 set but did for the LDBDS training set [AE] see simulation time settings above). Varying
307 the time-scale for the LDBDS model was necessary for analyzing real world data where
308 time-scales of outbreaks can vary considerably.

309 Our CNNs were implemented in Python 3.8.10 using keras v. 2.6.0 and
310 tensorflow-gpu v. 2.6.0. (Chollet; Abadi et al. 2016). [AE, R3] Convolutional Neural
311 Networks (CNNs) consist of one or more layers specifically intended for structural feature
312 extraction. CNNs utilize a filter, akin to a sliding window, that executes a mathematical
313 operation (convolution) on the input data. When dealing with structured data like the
314 CBLV+S matrix, multiple 1D filters slide across the matrix's columns, embedding each
315 scanned window into an N-dimensional vector representation. This architectural design
316 imparts CNNs with translation invariance, enabling them to recognize and learn repeating
317 patterns throughout the input space, regardless of their specific location. Stacking multiple
318 convolutional layers enables CNNs to decipher hierarchical structures within the data. See
319 Alzubaidi et al. (2021) and Khan et al. (2020) for reviews of the subject.

320 For each model, LIBDS and LDBDS, we designed and trained two CNN
321 architectures, one to predict epidemiological rate parameters and the other to predict the
322 outbreak location resulting in four total CNNs trained by two training data sets (LIBDS
323 and LDBDS). We used the mean-squared-error for the regression neural loss function in the
324 network trained to estimate epidemiological rates, and the categorical cross-entropy loss
325 function for the categorical network trained to estimate outbreak location. We assessed the
326 performance of the network by randomly selecting 5,000 samples for validation before each
327 round of training. We measured the mean absolute error and accuracy using the validation
328 sets. We used these measures to compare architectures and determine early stopping times
329 to avoid overfitting the model to the training data. We also added more simulations to the
330 training set until we could no longer detect an improvement in error statistics. After
331 comparing the performance of several networks, we found that the CNN described in SI
332 Figure S1 performed the best. In brief, the networks have three parallel sets of sequential
333 convolutional layers for the CBLV+S tensor and a parallel dense layer for the priors and
334 tree statistics. The three sets of convolution layers differed by dilation rate and stride
335 lengths. These three segments and the dense layer were concatenated and then fed into a

336 segment consisting of a sequential set of dense layers, each layer gradually narrowing to the
337 output size to either three or five for the rates and origin location networks, respectively,
338 for the LIBDS model, and seven and five for the seven rates and five locations, respectively,
339 for the LDBDS model.

340 All layers of the CNN used rectified linear unit (ReLU) activation functions. We
341 used the Adam optimizer algorithm for batch stochastic gradient descent (Kingma and Ba
342 2017) with batch size of 128. [AE, R2.3] We selected the number of epochs by monitoring
343 the mean absolute error and accuracy of the validation data set (which was not used in
344 training or testing) which suggested stopping after 15 epochs for the regression network
345 and ten epochs for the root location network would maximize accuracy/minimize error for
346 out-of-sample test data. [AE.8] The output layer activation for the network that predicted
347 the R_0 , δ and m parameters was linear with three nodes. For the output layer predicting
348 the outbreak location the activation function was softmax with five nodes for the five
349 locations. The input layer and all intermediate (latent) layers were the same for all four
350 networks, namely the CBLV+S tensor and the recovery rate, mean branch lengths, tree
351 height and number of tips in the tree. The LDBDS neural network was trained with
352 simulated trees where R_{0_i} varied among locations had output layer with seven nodes; five
353 for the each location's R_{0_i} and a node each for the sampling rate and the migration rate.
354 We tested networks with max-pooling layers between convolution layers as well as dropout
355 at several rates and found no improvement or a decrease in performance.

356 *Likelihood-based method of inference*

357 We compared the performance of our trained phylodynamic CNN to likelihood-based
358 Bayesian phylodynamic inferences. We specified LIBDS and LDBDS Bayesian models that
359 were identical to the LIBDS and LDBDS simulation models that we used to train our
360 CNNs. The most general phylodynamic model in the birth-death family applied to
361 epidemiological data is the state-dependent birth-death-sampling process (SDBDS;

362 (Kühnert et al. 2016; Scire et al. 2020)), where the state or type on which birth, death, and
 363 sampling parameters are dependent is the location in this context. The basic model used
 364 for experiments here is a phylogeographic model that is similar to the serially sampled
 365 birth-death process (Stadler 2010) where rates do not depend on location, which we refer
 366 to as the [authors:] LIBDS model. The death rate, μ , is equivalent to the recovery rate, γ ,
 367 in SIR models. Standard phylogenetic birth-death models assume the birth and death
 368 rates, λ and μ , are constant or time-homogeneous, while the SIR model's infection rate is
 369 proportional to β and S and varies with time as S changes. However, when the number of
 370 infected is small relative to susceptible people, as in the initial stages of an outbreak, the
 371 infection rate, β , is approximately constant and approximately equal to the birth rate λ ;

$$\lambda = \frac{\beta S}{N} \approx \beta \quad (8)$$

372 The joint prior distribution was set to the same model parameter distributions that
 373 were used to simulate the training and test sets of phylogenetic trees in the first section
 374 with γ treated as known and the proportion of extant lineages sampled, ρ , set to 0.01 as in
 375 the simulations. The likelihood was conditioned on the tree having extant samples (*i.e.* the
 376 simulation ran for the allotted time without being rejected). All simulated trees in this
 377 study had a stem branch and the outbreak origins were inferred for the parent node of the
 378 stem branch.

379 We used Markov chain Monte Carlo (MCMC) to simulate random sampling from
 380 the posterior distribution implemented in the TensorPhylo [authors:] plugin
 381 (<https://bitbucket.org/mrmay/tensorphylo/src/master/>) in RevBayes (Höhna et al. 2016).
 382 After a burnin phase, a single chain was run for 7,500 cycles with 4 proposals per cycle and
 383 at least 100 effective sample size (ESS) for all parameters. If the effective sample size (ESS)
 384 was less than 100, the MCMC was rerun with a higher number of cycles. We also analyzed

385 the coverage of the 5, 10, 25, 50, 75, 90, and 95% highest posterior density (HPD) intervals
386 to verify that our simulation model and inference model are the same and that the MCMC
387 simulated draws from the true posterior distribution. Bayesian phylogeographic analysis
388 recovered the true simulating parameters at the expected frequencies (Figure 2 C), thus
389 validating the simulations were working as expected and confirming that the MCMC was
390 accurately simulating draws from the true posterior distribution.

391 *Quantifying errors and error differences*

392 We measure the absolute percent error (APE) of the predictions from the CNN and the
393 mean posterior estimate (MPE) of the likelihood-based method. The formula for APE of a
394 prediction/estimate, y^{estimate} , of y^{truth} is

$$\text{APE} = \left| \frac{y^{\text{estimate}} - y^{\text{truth}}}{y^{\text{truth}}} \right| \times 100$$

395 The Bayesian alternative to significance testing is to analyze the posterior
396 distribution of parameter value differences between groups. In this framework, the
397 probability that a difference is greater than zero can be easily interpreted. We therefore
398 used Bayesian statistics to infer the median difference in error between the CNN and
399 likelihood-based methods and the increase in median error of each method when analyzing
400 misspecified data compared to when analyzing data simulated under the true inference
401 model.

402 We used Bayesian inference to quantify population error by performing three sets of
403 analyses: (1) inferred the population median APE under the true model (this will be the
404 reference group for analysis 3), (2) the effect of inference method — CNN or
405 likelihood-based (Bayesian) — on error by inferring the median difference between the
406 CNN estimate and the likelihood-based estimate, (3) the effect of misspecification on error

407 for each parameter by comparing the median error of estimates under misspecified
408 experiments and the reference group defined by analysis 1. See SI Figures S3 - S13 and SI
409 Table S1 for summaries and figures for all analyses for this section.

410 To infer these differences between groups we used the R package BEST (Meredith
411 and Kruschke). BEST assumes the data follow a t-distribution parameterized by a location
412 parameter, μ , a scale parameter, σ , and a shape parameter, ν , which they call the
413 "normality parameter" (*i.e.* if ν is large the distribution is more Normal). Because the
414 posterior distribution does not have a closed form, BEST uses Gibbs sampling to simulate
415 draws from the posterior distribution. 20,000 samples were drawn from the posterior
416 distribution for each BEST analysis. BEST uses automatic posterior predictive checks to
417 indicate that a model adequately describes the data distributions. Posterior predictive
418 checks indicate the BEST model adequately fits each data set analyzed below.

419 *Inferring the median APE.*— Before inferring differences between groups, we inferred the
420 population median APE for predictions of R_0 , δ , and m from test data simulated under the
421 inference model using the CNN and likelihood-based methods. Histograms of the sampled
422 log-transformed APE appears to be symmetric with heavy tails so we fit the log APE to
423 the BEST model. This implies that the sampled APE scores are drawn from a log-t
424 distribution. The log-t distribution has a mean of ∞ and median of e^μ , we therefore focus
425 our inference on estimating posterior intervals for the population median APE from the
426 sampled APE values for each parameter estimated by the CNN method and
427 likelihood-based method which we denote APE^{CNN} , and APE^{Like} respectively. The data
428 analyzed here and likelihood assumed by BEST is

$$y = \text{APE}^{\text{CNN}} \text{ or } \text{APE}^{\text{Like}}$$

$$\log y \mid \mu, \sigma, \nu \sim t_\nu(\mu, \sigma).$$

429 The priors were set to the vague priors that BEST provides by default,

$$\mu \sim \text{Normal}(\text{mean}(y), \text{sd}(y) \times 1000)$$

$$\sigma \sim \text{Uniform}(\text{sd}(y)/1000, \text{sd}(y) \times 1000)$$

$$\nu \sim \text{Exponential}(1/29) + 1.$$

430 95% highest posterior intervals (HPI) for the median APE, $\tilde{\mu}$, was estimated by the
431 following transformation of simulated draws from the posterior distribution

$$\tilde{\mu} = e^\mu.$$

432 In summary, the results we present are 95% HPI from the posterior distributions of
433 the median error, $\tilde{\mu}$.

434 *Inferring the relative accuracy of the CNN and likelihood-based method.*— To quantify the
435 difference in error between the CNN and the likelihood-based method, we fit the difference
436 in sampled APE scores, ΔAPE , between the CNN method and the likelihood-based
437 method to the BEST model. Histograms of ΔAPE appear symmetric with weak to strong
438 outliers making the BEST model a good candidate for inference from this data. The data
439 and likelihood are

$$\Delta y = \text{APE}^{\text{CNN}} - \text{APE}^{\text{Like}}$$

$$\Delta y \mid \mu, \sigma, \nu \sim t_\nu(\mu, \sigma)$$

440 We used the same default priors as above.

441 Because, Δy is not log-transformed, it is drawn from a t-distribution and the
 442 marginal posterior of the parameter μ is an estimate of the population mean, μ^d . Because
 443 the mean and the median are equivalent for a t-distribution, we again report the posterior
 444 distribution of the median difference, $\tilde{\mu}^d$ to simplify the results.

445 In summary, the results we present are 95% HPI from the posterior distribution of
 446 the median difference between the two methods, $\tilde{\mu}^d$.

447 When comparing CNN to the likelihood-based approach, positive values for $\tilde{\mu}^d$
 448 indicate the CNN is less accurate, and negative indicate the likelihood-based estimates less
 449 accurate. We emphasise that this quantity is the median difference in contrast to the
 450 difference in medians, $\Delta\tilde{\mu}$, reported in the next section.

451 *Inferring sensitivity to model misspecification.*— Finally, to quantify the overall sensitivity
 452 of each rate parameter to model misspecification under each inference method, we infer the
 453 difference in median APE, $\tilde{\mu}$ of predictions under a misspecified model relative to
 454 predictions under the true model. In other words we are inferring differences in medians
 455 between experiments. For example, to infer the sensitivity of the CNN’s inference of the
 456 sampling rate, δ , to phylogenetic error, we inferred the difference between the median APE
 457 of the CNN’s predictions for misspecified trees and the median APE of CNN predictions
 458 for true trees. The data is concatenated as below.

$$(y_1, y_2) = (\text{APE}^{\text{CNN}}, \text{APE}^{\text{CNN Ref}}) \text{ or}$$

$$(y_1, y_2) = (\text{APE}^{\text{Like}}, \text{APE}^{\text{Like Ref}})$$

459 We inferred the difference between group median APE scores, denoted $\Delta\tilde{\mu}$, by
 460 assuming that the model parameters conditioned on the observed APE from the two
 461 groups, y_1 and y_2 , follow a posterior distribution that is proportional to

$$P(y_1 | \mu_1, \sigma_1, \nu) P(y_2 | \mu_2, \sigma_2, \nu) P(\mu_1, \mu_2, \sigma_1, \sigma_2, \nu),$$

462 where $\log y_1$ and $\log y_2$ follow t distributions with means μ_1 and μ_2 and standard

463 deviations σ_1 and σ_2 , respectively while sharing a common normality parameter, ν .

464 The posterior sample of $\Delta\tilde{\mu}$ is obtained by transforming samples from the joint

465 marginal posterior distribution of μ_1 and μ_2 with the following equation,

$$\Delta\tilde{\mu} = e^{\mu_1} - e^{\mu_2}.$$

466 The two components of the likelihood are each t-distributed and share the ν

467 parameter which means we assume both samples are drawn from a similarly shaped

468 distribution (similarly heavy tails).

$$\log y_1 | \mu_1, \sigma_1, \nu \sim t_\nu(\mu_1, \sigma_1)$$

$$\log y_2 | \mu_2, \sigma_2, \nu \sim t_\nu(\mu_2, \sigma_2)$$

469 The prior distribution for the parameters of the model were set to the defaults for

470 BEST,

$$\begin{aligned}\mu_1 &\sim \text{Normal}(\text{mean}(\log y_1), \text{sd}(\log y_1) \times 1000) \\ \mu_2 &\sim \text{Normal}(\text{mean}(\log y_2), \text{sd}(\log y_2) \times 1000) \\ \sigma_1 &\sim \text{Uniform}(\text{sd}(\log y_1)/1000, \text{sd}(\log y_1) \times 1000) \\ \sigma_2 &\sim \text{Uniform}(\text{sd}(\log y_2)/1000, \text{sd}(\log y_2) \times 1000) \\ \nu &\sim \text{Exponential}(1/29) + 1\end{aligned}$$

⁴⁷¹ As before, interpretation of the posterior distribution of the difference in medians is
⁴⁷² straightforward: the more positive the difference in median APE from the misspecified
⁴⁷³ model test set and the median APE from the true model test set, the more sensitive the
⁴⁷⁴ parameter is to model misspecification in the experiment.

⁴⁷⁵ [EIC, R1.2]: CNN Uncertainty Quantification

⁴⁷⁶ We used conformalized quantile regression (CQR) to construct calibrated probability
⁴⁷⁷ intervals (CPI), ensuring accurate predictive coverage (Lei et al. 2018; Romano et al. 2019;
⁴⁷⁸ Sousa et al. 2022; Vovk et al. 2022; Angelopoulos et al. 2023). CQR is implemented in two
⁴⁷⁹ stages: first a network is trained to predict conditional quantiles, then a hold-out simulated
⁴⁸⁰ dataset is used to estimate bias adjustment terms to ensure correct coverage on future data
⁴⁸¹ i.e. 95% intervals contain the true value 95% of the time for test data.

⁴⁸² To implement quantile regression with a neural network and predict lower and upper
⁴⁸³ quantiles, we adjusted the general network architecture used for point estimates above to
⁴⁸⁴ have two outputs each with a mean pinball loss function instead of the mean squared error,

$$L_\tau(y, \hat{q}) = \frac{1}{N} \sum_i^N [(y_i - \hat{q}_i)\tau \mathbb{1}\{y_i \geq \hat{q}_i\} + (\hat{q}_i - y_i)(1 - \tau)\mathbb{1}\{y_i \leq \hat{q}_i\}] .$$

485 Here, y is the label or true parameter value (not a quantile) and \hat{q} is the trained neural
 486 network's prediction of a given quantile. τ is the quantile level and is equal to $1 - \alpha$, where
 487 α is the mis-coverage rate, or the probability the true value is not below the quantile. To
 488 estimate inner quantiles with miscoverage rate α , the lower quantile output was set to
 489 predict the $\alpha/2$ quantile for each rate parameter and the other layer to predict the $1 - \alpha/2$
 490 upper quantile (Steinwart and Christmann 2011) (SI figure S2). We refer to CNNs of this
 491 type as qCNN. Though often close, these inner quantiles are not guaranteed to have the
 492 correct coverage on test data sets (Figure 3) necessitating the calibration
 493 (conformalization) step (Romano et al. 2019).

494 To calibrate the predictions of quantile regression neural networks, CQR finds an
 495 adjustment term for each quantile through computing a non-comformity score, such as the
 496 distance of the predicted value from the predicted quantile. If the estimated quantile is
 497 well calibrated, then the same quantile of the scores in a calibration set will be zero. If the
 498 estimated quantile is, for example, too high then too high a proportion of the labels will
 499 fall below the estimated quantile and the empirical quantile, Q , of the nonconformity score
 500 $y - \hat{q}$ at $1 - \alpha/2$ will be negative. In other words it will over cover the calibration set. Q
 501 thus becomes the adjustment term for calibrating the qCNN's quantile estimate (equations
 502 9, and 10) by simply adding the term to the corresponding estimated quantile as shown in
 503 equation 11.

$$Q_{\text{lower}} \text{ s.t. } P(y - \hat{q}_{\text{lower}} < Q_{\text{lower}}) = \frac{\alpha}{2} \left(1 + \frac{1}{n}\right) \quad (9)$$

$$Q_{\text{upper}} \text{ s.t. } P(y - \hat{q}_{\text{upper}} < Q_{\text{upper}}) = \left(1 - \frac{\alpha}{2}\right) \left(1 + \frac{1}{n}\right) \quad (10)$$

504 $\text{CPI} = [\hat{q}_{\text{lower}} + Q_{\text{lower}}, \hat{q}_{\text{upper}} + Q_{\text{upper}}] \quad (11)$

505 Note that the quantiles of the score for finite sample sizes requires adjustment by $(1 + \frac{1}{n})$
 506 where n is the number of samples in the calibration set (Romano et al. 2019).

507 We simulated 108,559 more datasets (trees) to estimate the calibration amounts for
508 the upper and lower qCNN-estimated quantiles. After calibration through
509 conformalization, we clipped intervals to the prior boundary for intervals that extended
510 beyond the prior distribution's range. To examine the consistency of quantile regression for
511 neural networks trained on different quantiles we trained seven different quantile networks
512 to predict the same quantiles used for validating our Bayesian analysis and simulation
513 model; {0.05, 0.25, 0.5, 0.75, 0.9, 0.95}. We checked the coverage of these adjusted CPIs on
514 another simulated test dataset of 5,000 trees.

515 *Real Data*

516 We compared the inferences of a LDBDS simulation trained neural network to that of a
517 phylodynamic study of the first COVID wave in Europe (Nadeau et al. 2021). These
518 authors analyzed a phylogenetic tree of viruses sampled in Europe and Hubei, China using
519 a location-dependent birth-death-sampling model in a Bayesian framework using priors
520 informed by myriad other sources of information. We simulated a new training set of trees
521 under an LDBDS model where $R_{0,i}$ depends on the geographic location, and the sampling
522 process only consists of serial sampling and no sampling of extant infected individuals.**[AE,**
R1.2] We estimated 95% CPIs for model parameters with a simulated calibration dataset of
524 101,219 trees using CQR as above and confirmed accurate coverages with another dataset
525 of 5,000 trees.

526 We then analyzed the whole tree from Fig. 1 in (Nadeau et al. 2021) as well as the
527 European clade which Nadeau et al. (2021) labeled as A2 in the same figure. We note that
528 our simulating model is not identical to the inference model used in (Nadeau et al. 2021).
529 We model migration with a single parameter with symmetrical migration rates among
530 locations and all locations having the same sampling rate. Nadeau and colleagues
531 parameterize the migration process with asymmetric pairwise migration rates and assume
532 location-specific sampling rates. We also do not include the information the authors used

533 to inform their priors as that requires an extra level of simulation and training on top of
534 simulations done here, and is thus beyond the scope of this study.

535 The time tree from (Nadeau et al. 2021) was downloaded from GitHub
536 (<https://github.com/SarahNadeau/cov-europe-bdmm>). The recovery rate assumed in
537 (Nadeau et al. 2021) was 0.1 days^{-1} which was set to 0.05 to bring the recovery rate to
538 within the range of simulating values used to train the CNN. Consequently, the branch
539 lengths of the tree were then scaled by 2. The number of tips, tree height, and average
540 branch lengths were measured from the rescaled trees and fed into the network. The full
541 tree and A2 clade were then analyzed using the LDBD CNN and compared to the posterior
542 distributions from (Nadeau et al. 2021).

543 *Hardware used*

544 Simulations were run on a 16 core Intel(R) Xeon(R) Platinum 8175M CPU @ 2.50GHz.
545 For each simulation, an XML file with random parameter settings was generated using
546 custom scripts. These XML files were the inputs for MASTER which was run in the
547 BEAST2 platform. Neural network training and testing and predictions were conducted on
548 an 8 core Intel(R) Core(TM) i7-7820HQ CPU @ 2.90GHz laptop [R2.7] with a NVIDIA
549 Quadro M1200 GPU for training.

550 RESULTS

551 *Comparing deep learning to likelihood*

552 Our first goal in this study was to train a CNN that produced phylodynamic parameter
553 point estimates that were as accurate as likelihood-based Bayesian posterior mean
554 estimates under the true model. This will serve as a reference for quantifying level of
555 sensitivity in our misspecification experiments. [authors:] Using viral phylogenies like those

556 typically estimated from serially sampled DNA sequences, we focused on estimating
557 important epidemiological parameters – the reproduction number, R_0 , the sampling rate, δ ,
558 the migration rate, m , and the outbreak origin.

559 Our CNN produced estimates that are as accurate as the mean posterior estimates
560 (MPE) under the true simulating model. We compared the absolute percent error (APE)
561 of the network predictions to the APE of the MPE of the Bayesian location-independent
562 birth-death-sampling (LIBDS) model (Figure 2). The APE is straight-forward to interpret,
563 e.g. an APE of < 10 means the estimate is within 10 percentage points (pps) of the true
564 value. For the three epidemiological rate parameters, R_0 , δ and m , both methods made
565 very similar predictions for the 100 time tree test set (Figure 2 panel A). The two methods
566 appear to produce estimates that are more similar to each other than to the ground truth
567 labels (compare bottom row scatter plots in orange to the blue and red scatter plots in
568 panel A). Fig. 2 panel B shows that the inferred median difference in APE, $\tilde{\mu}^d$, between
569 the method's estimates for the three parameters is close to zero ($|\tilde{\mu}^d|$ 95% highest
570 posterior interval (HPI) is < 4 pps; SI Table S1; SI Figure S3).

571 [EIC, R1.2] We also compared the performance of uncertainty quantification using
572 quantile-CNN-based conformalized quantile regression (CQR; Romano et al. 2019) to that
573 of Bayesian HPIs for each of the experiments. We trained seven qCNNs to predict
574 inner-quantiles at seven different levels to compare with the Bayesian HPIs; $\tau = \{0.05, 0.1,$
575 $0.25, 0.5, 0.75, 0.9, 0.95\}$. We then used another simulated dataset to calibrate predicted
576 intervals which we refer to as calibrated probability intervals (CPIs) which theoretically
577 have correct coverage properties (Romano et al. 2019) like the HPIs. For the test dataset of
578 138 trees, the CPIs had coverages that matched well with expectations to a comparable
579 degree to the Bayesian highest posterior density intervals (HPI) (Figure 2 panel C) though
580 more variable. To further confirm that our CQR procedure was adequately calibrating the
581 qCNN estimates, we confirmed correct coverages of CPIs for a much larger dataset with
582 5,000 trees (Figure 3). On average, the widths of CPIs in the set of 138 trees shown in

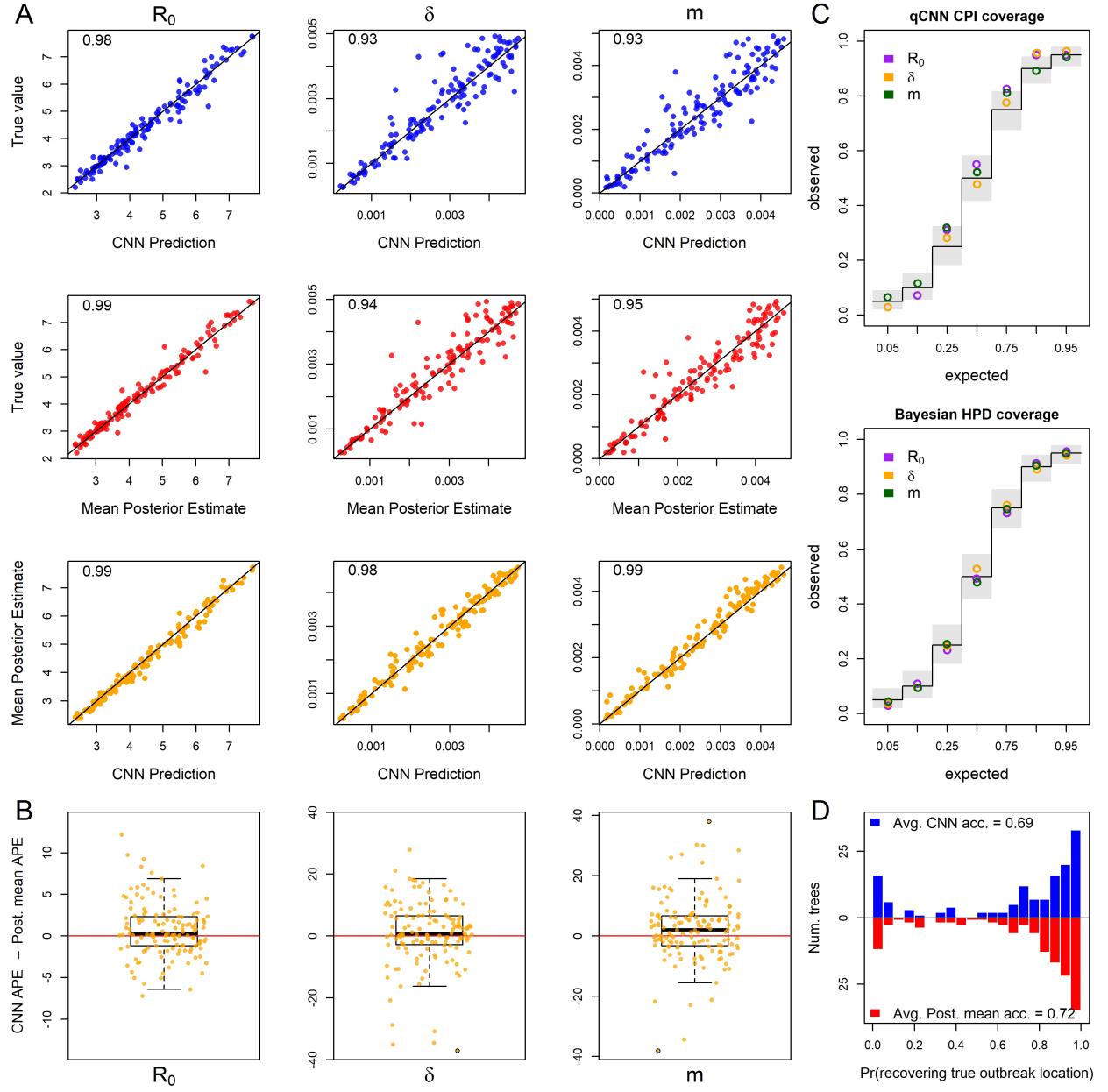


Figure 2: Inference under the true simulating model. (A) Scatterplot of CNN predictions and posterior mean estimates from Bayesian analyses against the true values (top two rows in blue and red respectively) of the basic reproduction number, R_0 , the sampling rate, δ , and the migration rate, m for 138 test trees. [AE, R1.2] In the upper-left corners of the scatter plots are the correlations of the plotted data. The bottom row in orange shows scatter plots of the CNN estimates against the posterior mean estimates for the same trees. (B) The difference in the absolute percent error (APE) of estimates for the two inference methods. Boxes show the inner 50% quantile of the data while whiskers extend 1.5 IQR. Dots with black circles were truncated to 2× the length of whiskers for visualization purposes. [AE, R1.2] (C) Coverage plots show the expected frequency of coverage for each of the categories and the observed frequencies (black steps and colored circle respectively). Gray boxes are the expected 95% confidence intervals at each of the expected coverage values which follows a Beta($(n+1)q, n-(n+1)q+1$) distribution. (D) Histograms of the probabilities of inferring the correct outbreak origin location.

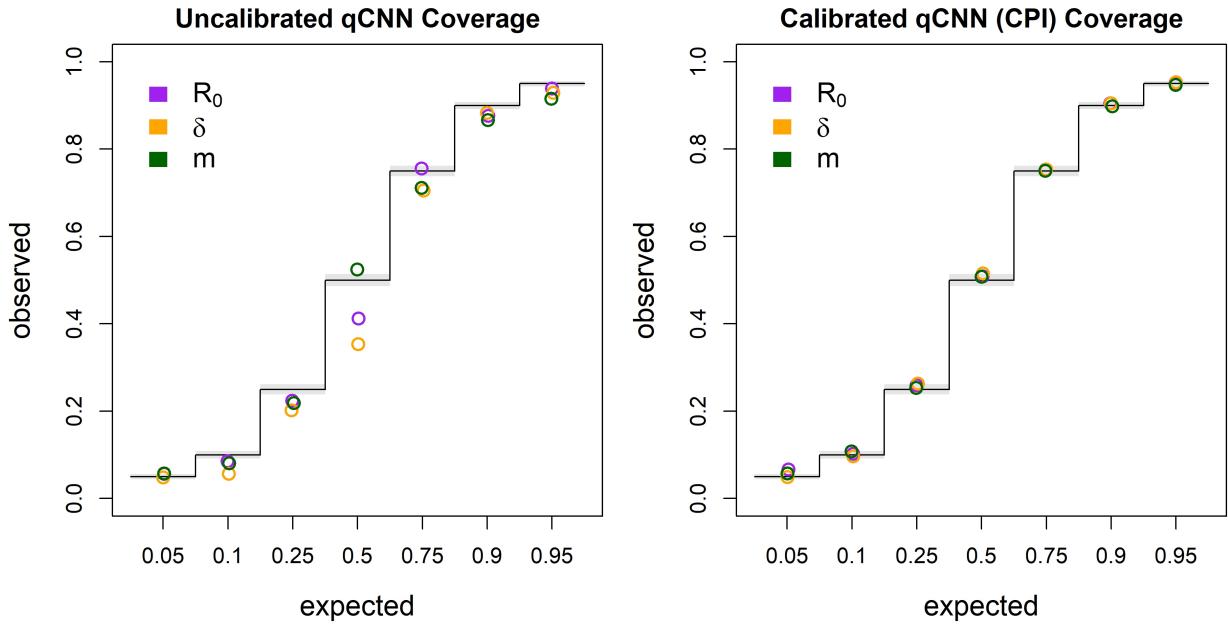


Figure 3: Coverage of uncalibrated qCNN quantile predictions (left) and calibrated qCNN which produce “calibrated probability intervals” (CPI) on the right. The observed coverage of 5,000 samples tested at seven different predicted coverage levels (labeled horizontal). See Figure 2 C for more details on coverage plots.

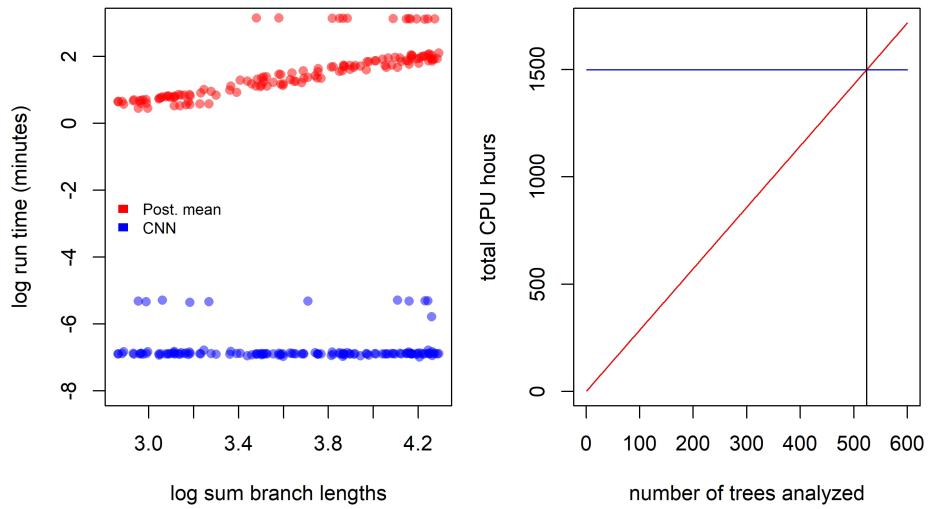


Figure 4: Left: Estimates of time to complete analysis of each of 138 trees relative to tree size. Right: The number of trees (524; gray vertical line) needed to analyze for total analysis time of Bayesian method (red line) to equal that of the entire simulation and CNN training and inference pipeline (blue line).

583 figure (2) was about 20 - 40% wider than that of the corresponding HPI and Jaccard
584 similarity index ranging from 0.66 to 0.75 suggesting a high degree of overlap between the
585 intervals (SI Figure S4 and SI Table S2). These results indicate the probability level of the
586 CPI, e.g. 95%, can be safely interpreted as the probability a parameter falls within the
587 CPI. The wider intervals suggest the basic CQR method employed here is somewhat less
588 precise and thus more conservative than the Bayesian method.

589 Our trained CNN provides nearly instantaneous estimates of model parameters.

590 While the run time of the likelihood approach employed in this study scales linearly with
591 the size of the tree, the neural network has virtually constant run times that are more than
592 three orders of magnitude faster. Because simulation-trained neural networks have a
593 one-time cost of simulating the training data set and then training the neural network,
594 these methods are often called amortized-approximators (Bürkner et al. 2022). This means
595 the time savings aren't recouped until a certain number of trees have been analyzed. For
596 example, here over 524 trees would need to be analyzed to realize the cost savings of
597 simulating data and training our neural network (Figure 4). This illustrates the importance
598 of simulation optimization and generality for likelihood-free approaches to inference.

599 *Comparing [authors] sensitivity to model misspecification*

600 To test the relative sensitivity of CNN estimates and the likelihood-based MPE to model
601 misspecification, we simulated several test data sets under different, more complex
602 epidemic scenarios and compared the decrease in accuracy (increase in APE).

603 Our first model misspecification experiment tested performance when assuming all
604 locations had the same R_0 when, in fact, each location had different R_{0i} values. The
605 median APE for all three parameters increased to varying degrees (SI Fig. S5 Panel A)
606 compared to the median APE measured in Fig. S3. We found that both methods
607 converged on similar biased estimates for R_0 . In both the CNN and Bayesian method,
608 estimates of δ were relatively robust to misspecifying R_0 . In contrast, the migration rate

showed much more sensitivity to this model violation in both methods with both methods also converging on similarly biased estimates (Figure 5 A). The median difference in error between the two methods is close to zero for all rate parameters ($|\tilde{\mu}^d| 95\% \text{ HPI} < 6 \text{ ppts}$; SI Table S1) (SI Figure S5 Panel B). [EIC, R1.2] For both methods of uncertainty quantification the coverage declined by similar amounts for all three parameters with δ showing little to no sensitivity to R_0 misspecification (Figure 5 panel C and SI Table S2). The patterns of coverage are also somewhat less regular across the qCNN quantiles than the HPIs for the migration rate parameter likely due in part to the fact that each inner quantile qCNN was trained independently and thus have independent errors. The relative interval widths and Jaccard similarity indexes did not change appreciably from predictions under the true model (SI Figure S4 and SI Table S2). Our CNN appears to be slightly more sensitive than the Bayesian approach when predicting the outbreak location. Nevertheless, their distributions are quite similar (Figure 5 Panel C).

Next, we measured method sensitivity when the sampling process of the test trees violates assumptions in the inference model. In this set, each location had a unique and independent sampling rate, δ , rather than a single δ shared among locations. The median APE only increased for δ and m (SI Figure S7 Panel A). As expected, estimates of δ were highly biased for both methods (Figure 6 panel A). Panel A also shows that R_0 is virtually insensitive to sampling model misspecification, but that migration rate, again, is highly sensitive in both the CNN and likelihood method. The median difference in error between the two methods is close to zero for all the rate parameters ($|\tilde{\mu}^d| 95\% \text{ HPI} < 5 \text{ ppts}$; SI Table S1, SI Figure S7) (Figure 6 panel B). [EIC, R1.2] For both methods coverage declined for δ and m , while R_0 showed little to no sensitivity to δ misspecification (Figure 6 panel C and SI Table S2). The relative widths and degree of overlap was again similar to the experiments above (SI Figure S8, SI Table S2). We again also see greater irregularity among CPI levels in coverage, notably δ at inner-quantile level 0.9. The location of outbreak prediction is also somewhat sensitive in both methods, with the CNN showing a

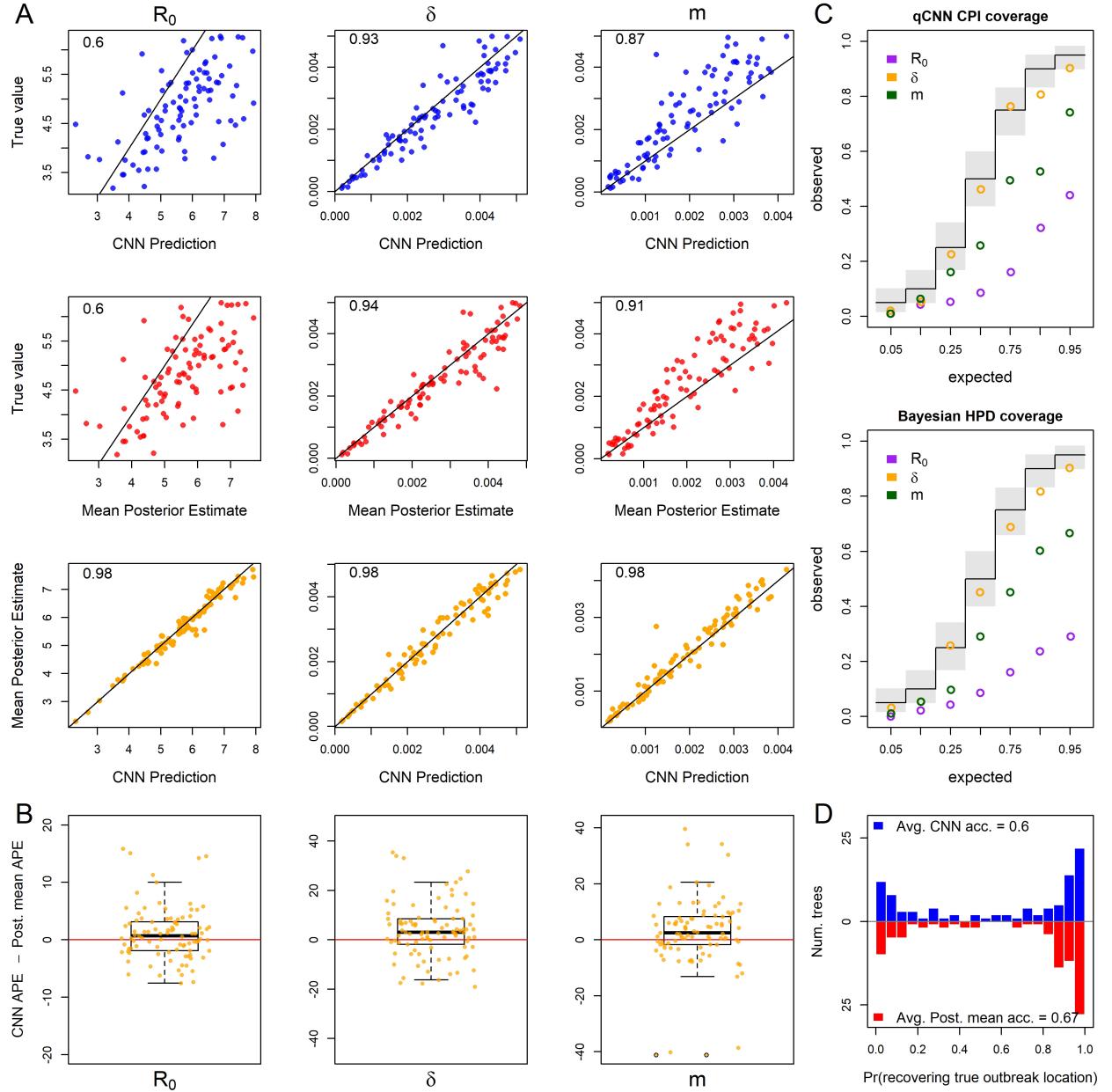


Figure 5: For 93 test trees where the R_0 parameter was misspecified: the simulating model for the test data specified 5 unique R_0 s among the five locations while the inference methods assumed one R_0 shared among locations. Because of this, the estimates for R_0 are plotted against mean of the five true R_0 values. See Figure 2 for general details about plots.

636 slightly larger mean difference, but the overall distribution of accuracy of all the test trees
 637 again is similar (Figure 6 panel C).

638 To explore sensitivity to migration model underspecification, we simulated a test set
 639 where the migration rates between locations is free to vary rather than being the same

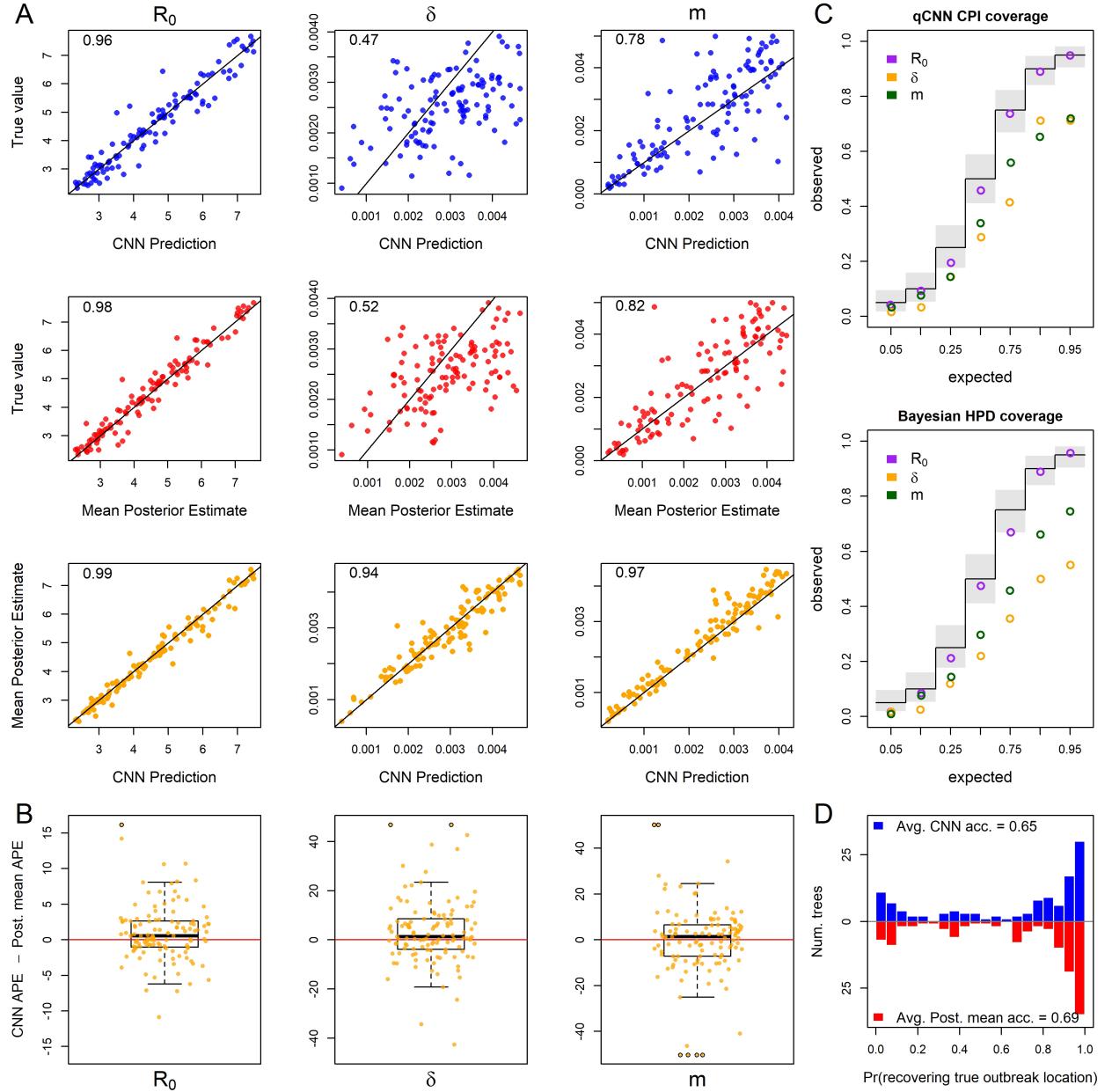


Figure 6: For 118 test trees where the sampling rate parameter was misspecified: the simulating model for the test data specified 5 unique sampling rates among the five locations while the inference methods assumed one sampling rate shared among locations. The estimates of δ are plotted against the mean true values of δ . See Figure 2 for general details about plots.

among locations as in the inference model. This implies $5!$ unique location-pairs and thus unique migration rates in the test data set. Results show that for both methods the parameters R_0 and δ are highly robust to this simplification (SI Fig. S9 Panel A). Though

643 estimates of a single migration rate had a high degree of error compared to a single pair of
644 locations' migration rates (Figure 7 panel A), the two methods still had similar estimates
645 with the difference in APE centered near zero (Figure 7 panel B). The inferred median
646 difference in APE was close to zero ($|\tilde{\mu}^d|$ 95% HPI < 3 ppt; SI Table S1; SI Figure S9
647 Panel B). [EIC, R1.2] For both methods the coverage only declined significantly for the
648 migration rate and the decrease was again similar in magnitude across quantiles (Figure 7
649 panel C and SI Table S2). Again, relative widths and degree of overlap of CPI and HPI
650 was similar to previous experiments (SI Figure S10, SI Table S2) There was a slight but
651 similar decrease in accuracy in predicting the outbreak location for both methods (Figure 7
652 panel C).

653 When testing the sensitivity of the two methods to arbitrary groupings of locations,
654 we found that both methods showed equal sensitivity to the same parameters (Fig. 8
655 Panels A and B). In particular, the migration rate showed a modest increase in median
656 APE and R_0 and sample rate showed virtually no sensitivity to arbitrary grouping of
657 locations (SI Figure S11 Panel A). The inferred median difference between method APE's
658 was again close to zero ($|\tilde{\mu}^d|$ 95% HPI < 4 ppt; SI Table S1; SI Figure S11 Panel B).
659 [EIC, R1.2] For both methods the coverage declined modestly only for the migration rate
660 (Figure 5 panel C and SI Table S2). Relative widths and interval overlap showed virtually
661 no change (SI Figure S12 and SI Table S1). These results suggest that for at least the
662 exponential phase of outbreaks where rate parameters do not vary among locations, these
663 models have a fair amount of robustness to the decisions leading to geographical division of
664 continuous space into discrete space. The outbreak location showed higher accuracy in
665 both methods due to the fact that the test data was no longer a flat distribution; the 6
666 combined locations should contain 60% of the outbreak locations (Figure 8 panel C).

667 Finally, we explored the relative sensitivity of our CNN to amounts of phylogenetic
668 error that are present in typical phylogeographic analyses. Our simulated phylogenetic
669 error produced trees with average [R1.13] Jaccard similarity indexes between the inferred

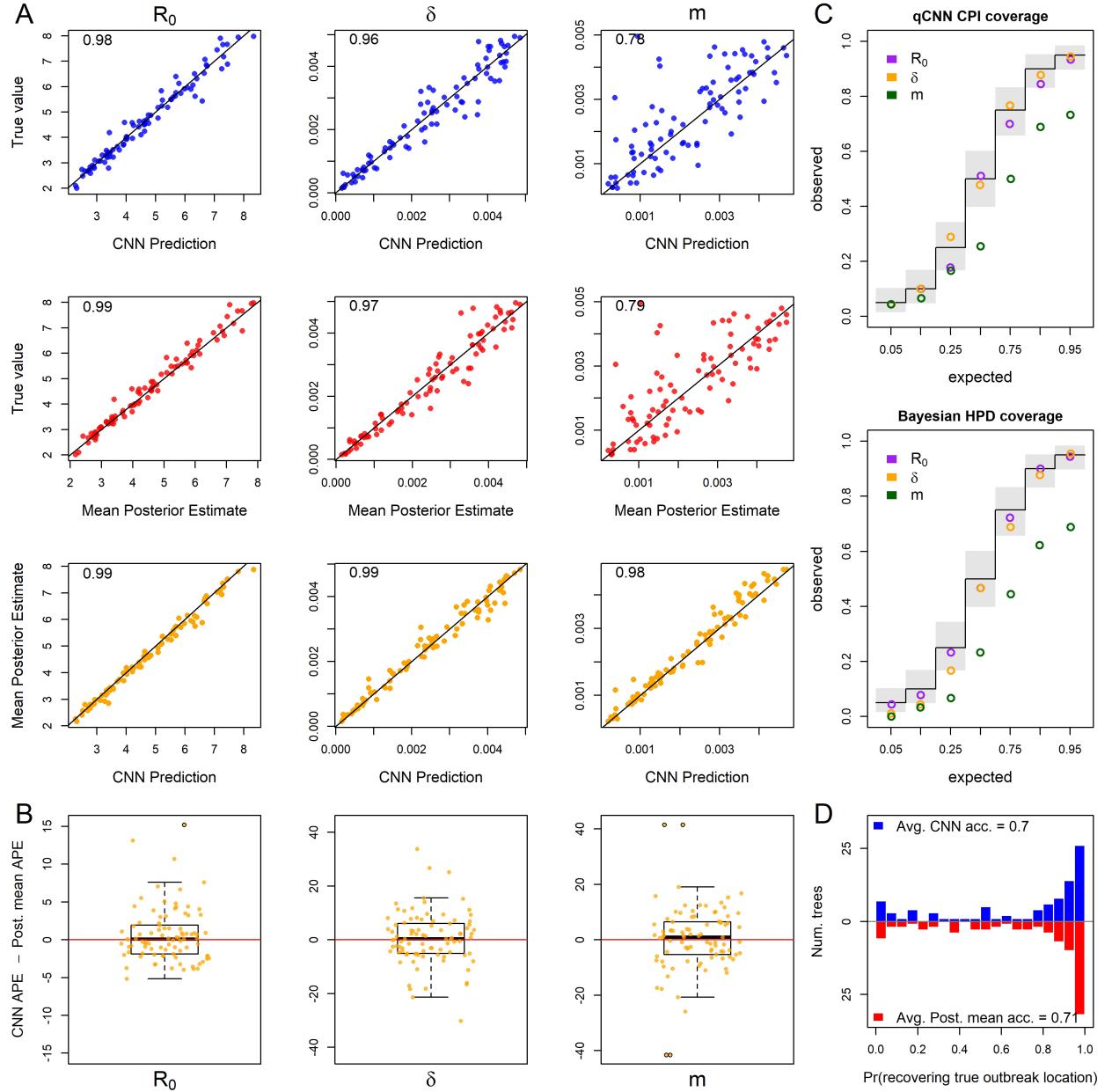


Figure 7: For 90 test trees where the migration rate parameter was misspecified: the simulating model for the test data specified 5! (120) unique migration rates among the unique pairs of the five locations while the inference methods assumed all migration rates were equal. The inferred migration rate is plotted against the mean pairwise migration rates of test data set. See Figure 2 for general details about plots.

670 tree and the true tree of about 0.5 with 95% of simulated trees having distances within
 671 0.36 and 0.72. We again compared inferences derived from the true tree and the tree with
 672 errors using the CNN and the Bayesian LIBDS methods. Results show that migration rate

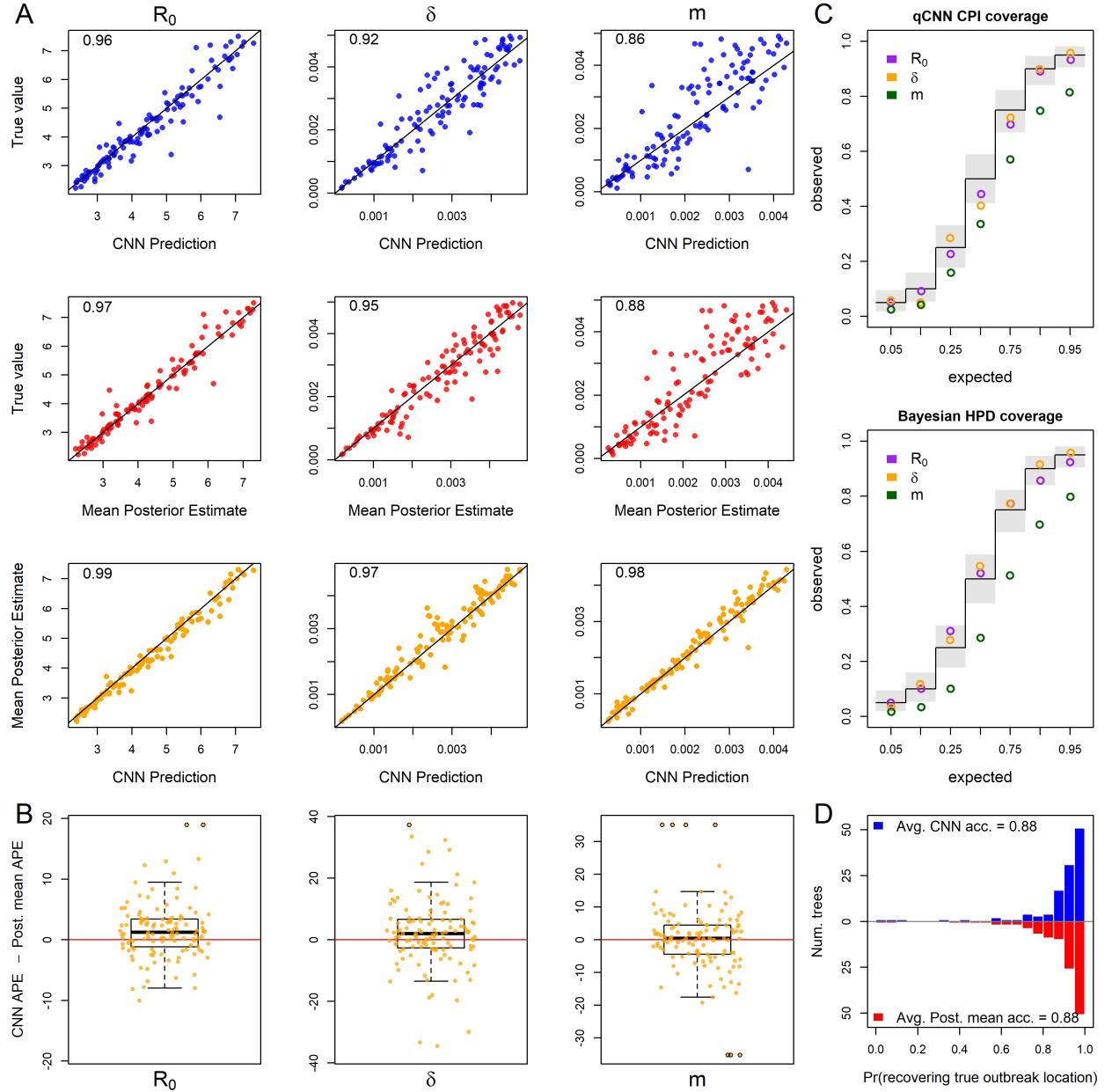


Figure 8: For 101 test trees where the number of locations was misspecified: the simulating model for the test data specified an outbreak among 10 locations with 6 locations subsequently combined into a single location while the inference methods assumed 5 locations with no arbitrary combining of locations. See Figure 2 for general details about plots.

673 was minimally affected but R_0 and δ were to a some degree sensitive to phylogenetic error
 674 (Figure 9 panel A; SI Figure S13 Panel A), with both methods again showing similar
 675 degrees of sensitivity (Figure 9 panel B). The inferred median difference was, yet again,
 676 small ($|\tilde{\mu}^d|$ 95% HPI < 6 ppts. SI Table S1, SI Figure S13 Panel B). [EIC, R1.2]

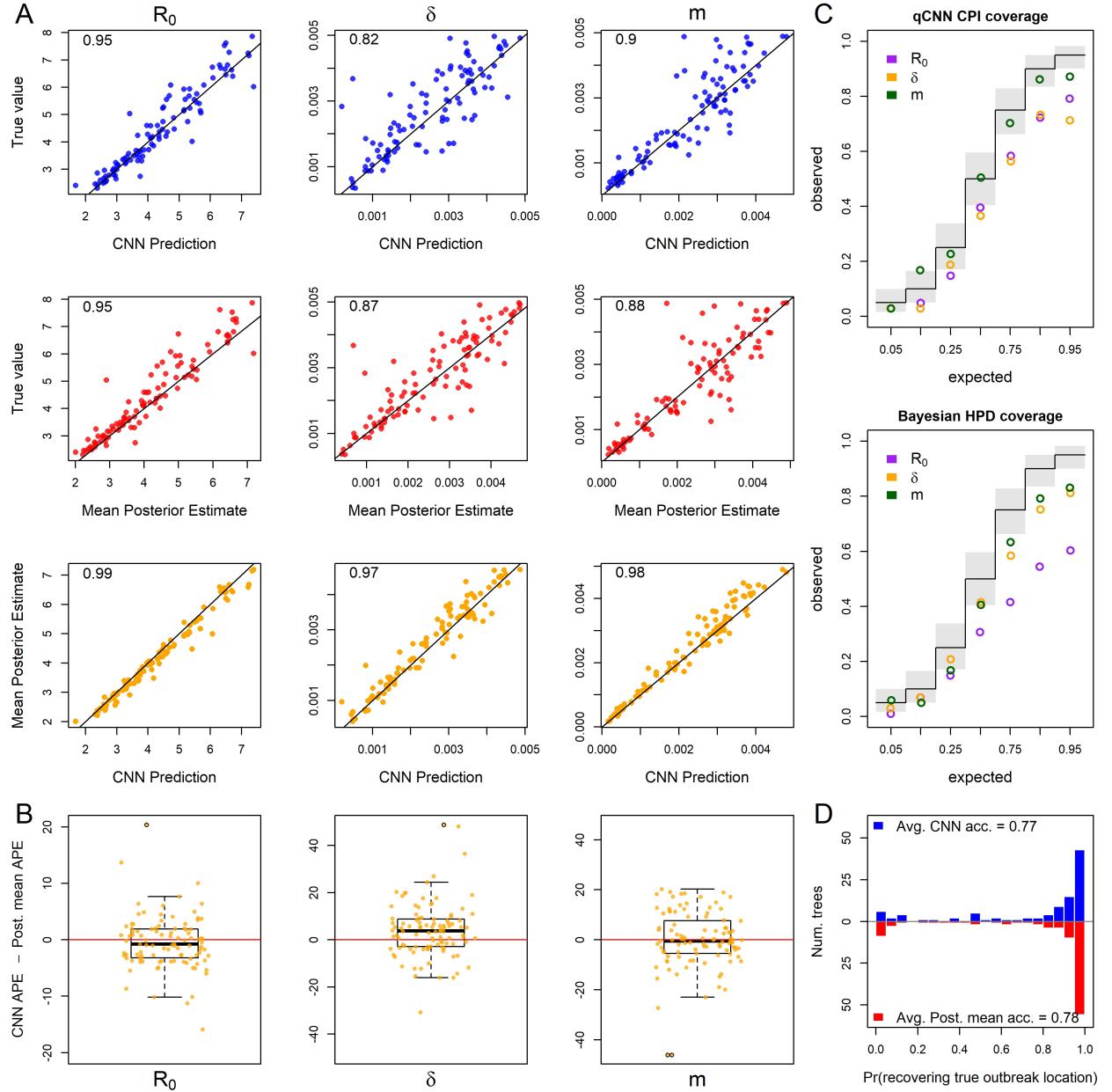


Figure 9: For 118 test trees where the time tree was misspecified: the true tree from the simulated test set was replaced with an inferred tree from simulated DNA alignments under the true tree. See Figure 2 for general details about plots.

677 Cov coverages of δ declined for both methods in a similar way across quantiles. Again the 90%
 678 inner quantile showed some inconsistency with its neighboring quantiles. In this case its
 679 coverage for δ was slightly higher than the 95th inner quantile. The CPIs for R_0 appear
 680 much less sensitive (Figure 9 panel C and SI Table S2). Although the relative widths of the

681 CPIs and HPIS was similar to previous experiments, the degree of overlap decreased
682 somewhat by about 5 - 10% (SI Figure S14 and SI Table S2). One difference between this
683 experiment and the others, is that trees are data instead of model parameters. It is
684 interesting that the point estimates from the two methods show similar biases while the
685 coverages seem to depart somewhat. Inference of the origin location, were very similar for
686 both methods (Fig. 9 Panel C).

687 *Analysis of SARS CoV-2 tree*

688 We next compared our likelihood-free method to a recent study investigating the
689 phylodynamics of the first wave of the SARS CoV-2 pandemic in Europe (Nadeau et al.
690 2021). Despite simulating the migration and the sampling processes differently from
691 Nadeau et al. (2021), our CNN produces similar estimates for the location-specific R_0 and
692 the origin of the A2 clade (Figure 10). Whether the full tree or just the A2 clade is fed into
693 the network, the predicted R_0 for each location was not far from the posterior estimates of
694 Nadeau et al. (2021). [AE, R1.2] For the most part the R_0 95% CPI for each location
695 overlaps to a high degree with the 95% HPI and is roughly 1.5 times wider indicating that
696 our CNN estimates are relatively conservative. For Hubei the interval width of the a2 clade
697 is much wider than the estimate using the whole tree. This is not surprising because there
698 are no samples from Hubei in the a2 clade. We also obtained estimates for a single
699 sampling rate and a single migration rate from our CNN and CPIs from our calibrated
700 qCNN. Among the five location-specific estimates of the sampling proportion and the
701 migration proportion from Nadeau et al. (2021), our CNN's point estimates and interval
702 estimates fall well within the their combined ranges.

703 The spillover location prediction CNN produced probability estimates of the A2
704 clade ancestral location the mostly agreed with that of Nadeau and colleagues (Figure 10,
705 right histograms). The only significant discrepancy in the European origin prediction is
706 that Nadeau and colleague's analysis suggests a much higher probability that the most

707 recent common ancestor of the A2 clade was in Hubei than our CNN predicts. This is
 708 likely because our CNN only used the A2 clade to predict A2 origins which has no Hubei
 709 samples to infer the origin of the A2 clade while Nadeau et al. (2021) used the whole tree.
 710 Notwithstanding this difference, among European locations, both methods predict
 711 Germany is the most likely location of the most recent common ancestor followed by Italy.

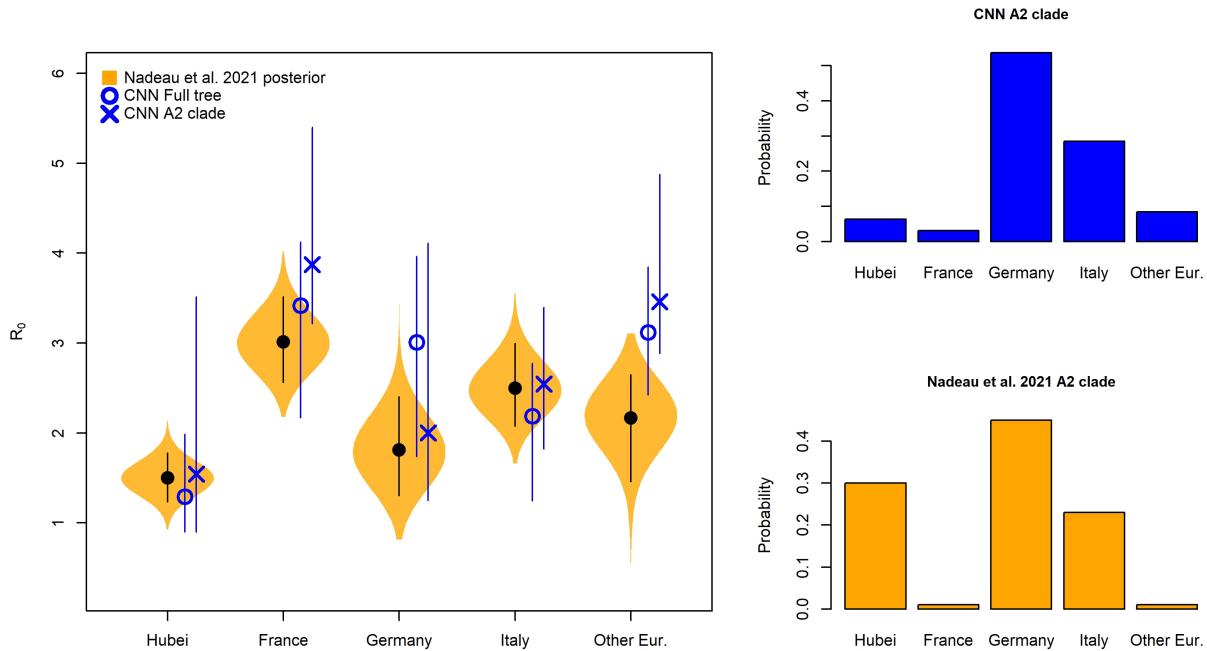


Figure 10: Location-dependent birth-death-sampling model (LDBDS) CNN comparison to (Nadeau et al. 2021) inference. Left violin plots show the posterior distributions of R_0 for each location in Europe as well as Hubei, China (orange). The block dot and line within each violin plot shows the posterior mean and 95% HPI respectively. The blue X and O marks the LDBDS CNN prediction from analyzing the full tree and the A2 (European) clade respectively. Vertical blue lines give the 95% calibrated probability interval (CPI) for the CNN estimates of R_0 . Right barplots show the LDBDS CNN prediction (blue) and posterior inference (orange) from (Nadeau et al. 2021) of the ancestral location of the A2 (European) clade (see Figure 1 (Nadeau et al. 2021)).

712

DISCUSSION AND CONCLUSIONS

713 Inference models are necessarily simplified approximations of the real world. Both
 714 simulation-trained neural networks and likelihood-based inference approaches suffer from

715 model under-specification and/or misspecification. When comparing inference methods it is
716 important to assess the sensitivity of model inference to simplifying assumptions. In this
717 study we show that newer deep learning approaches and standard Bayesian approaches
718 behave and misbehave in similar ways under a panel of phylodynamic estimation tasks
719 where the inference model is correct as well as when it is misspecified.

720 By extending new approaches to encode phylogenetic trees in a compact data
721 structure (Voznica et al. 2022; Lambert et al. 2022), we have developed the first application
722 of phylodynamic deep learning applied to phylogeography with serial sampling. Our
723 approach is similar to that of Lambert et al. (2022) in which they analyzed a binary SSE
724 model with exclusively extant sampling. By training a neural network on phylogenetic trees
725 generated by simulated epidemics, we were able to accurately estimate key epidemiological
726 parameters, such as the reproduction number and migration rate, in a fraction of the time
727 it would take with likelihood-based methods. Like Voznica et al. (2022) and Lambert et al.
728 (2022), we found that CNN estimators perform as well or nearly as well as likelihood-based
729 estimators under conditions where the inference model is correctly specified to match the
730 simulation model. The success of these separate applications of deep learning to different
731 phylodynamic problems is a testament to the versatility of the CBLV encoding of trees.

732 We compared the sensitivity of deep learning and likelihood-based inference to
733 model misspecification. Because deep-learning methods of phylogenetic and phylodynamic
734 inference are new, few studies compare how simulation-trained deep learning methods fail
735 in comparison to likelihood methods in this way (Flagel et al. 2019). We assume that when
736 the inference model is correctly specified to match the simulation model, the trained CNN
737 will, at best, produce noisy approximations of likelihood-based parameter estimates. In
738 reality, issues related to training data set size, learning efficiency, and network overfitting
739 may cause our CNN-based estimates to contain excess variance or bias when compared to
740 Bayesian likelihood-based estimators. Our results from five model misspecification
741 experiments show that both methods of inference perform similarly when the simulating

742 model and the inference model assumptions do not perfectly match. These similarities
743 exist not only in aggregate, when comparing method performance across datasets, but also
744 when comparing performance for each individual dataset. This suggests that the CNN and
745 likelihood methods are truly estimating parameters using isomorphic criteria, despite the
746 fact that CNN heuristically learns these criteria through data patterns, while likelihood
747 precisely and mathematically defines these criteria through the model definition itself.

748 Results of comparative sensitivity experiments like this are important because if
749 likelihood-free methods using deep neural networks can easily be trained to yield estimates
750 that are as robust to model misspecification as likelihood-based methods, then analysis of a
751 large space of more complex outbreak scenarios for which tractable likelihood functions are
752 not available can be developed and applied to real world data. Additionally, sufficiently
753 realistic, pre-trained neural networks can yield nearly instantaneous inferences from data in
754 real time to inform analysts and policy makers.

755 We also tested location-dependent SIR simulation trained neural network against a
756 previous publication fitting a similar model – location-dependent birth-death-sampling
757 (LDBDS) model – on real-world data using a Bayesian method. Our CNN predicted
758 location-specific R_{0_i} and outbreak origin in Europe were similar to that inferred in (Nadeau
759 et al. 2021). This result and our model misspecification experiments suggest that
760 simulation-trained deep neural networks trained on phylogenetic trees can find patterns in
761 the training data that generalize well beyond the training data set.

762 Our study extends the results of Voznica et al. (2022) and Lambert et al. (2022) in
763 several important ways. Our work showed that the new compact bijective ladderized vector
764 encoding of phylogenetic trees can easily be extended with one-hot encoding to include
765 metadata about viral samples. Using this strategy, we trained a neural network to not only
766 predict important epidemiological parameters such as R_{0_i} and the sampling rate, but also
767 geographic parameters such as the migration rate and the location of outbreak origination
768 or spillover. We anticipate that more diverse and complex metadata can be incorporated to

769 train neural networks to make predictions about many important aspects of
770 epidemiological spread such as the relative roles of different demographic groups and the
771 overlap of different species' ranges.

772 This approach can be readily applied to numerous compartment models used to
773 describe the spread of different pathogens among different species, locations, and
774 demographic groups, e.g. SEIR, SIRS, SIS, etc. (Ponciano and Capistrán 2011; Volz and
775 Siveroni 2018; Bjørnstad et al. 2020; Chang et al. 2020; O'Dea and Drake 2021) as well as
776 modeling super-spreader dynamics as in (Voznica et al. 2022). [Authors] Here we focused
777 on one phase of outbreaks (the exponential phase), but there are many other scenarios to
778 be investigated, such as when the stage of an epidemic differs among locations (e.g.
779 exponential, peaked, declining). The link between the underlying population dynamics
780 from which viral genomes are sampled and inferred phylogenetic trees can now easily be
781 interrogated through deep learning. [R1.3] More complex models will require larger trees to
782 infer model parameters. In this study we explored trees that contained fewer than 500 tips,
783 but anticipate that larger trees will demonstrate even greater speed advantages of neural
784 networks over likelihood-based methods either through subsampling regimes (Voznica et al.
785 2022) or by including larger trees in training datasets.

786 With fast, likelihood-free inference afforded by deep learning, the technical
787 challenges shift from exploring models for which tractable likelihood functions can be
788 derived towards models that produce realistic empirical data patterns, have parameters
789 that control variation of those patterns, and are efficient enough to generate large training
790 data sets. A growing number of advanced simulators are rapidly expanding the possibilities
791 for deep learning in phylogenetics. For example, FAVITES (Moshiri et al. 2019) is a
792 simulator of disease spread through large contact networks that tracks transmission trees
793 and simulates sequence evolution. Gen3sis, MASTER, SLiM, and VGsim are flexible
794 simulation engines for generating complex ecological, evolutionary, and disease
795 transmission simulations (Hagen et al. 2021; Vaughan and Drummond 2013; ?; Haller and

796 Messer 2019; Overcast et al. 2021). Continued advances in epidemic simulation speed and
797 flexibility will be essential for likelihood-free methods to push the boundaries of epidemic
798 modeling sophistication and usefulness.

799 There are several avenues of development still needed to realize the potential of
800 likelihood-free inference in phylogeography using deep learning. The current setup is ideal
801 for simulation experiments, but it is more difficult to ensure that the optimal parameter
802 values for empirical data sets are within the range of training data parameters.
803 Standardizing input tree height, geographical distance, and other parameters help make
804 training data more universally applicable. Simulation-trained neural networks are often
805 called amortized methods (Bürkner et al. 2022; Schmitt et al. 2022) because the cost of
806 inference is front-loaded, *i.e.* it takes time to simulate a training set and train a neural
807 network. The total cost in time per phylogenetic tree amortizes as the number of trees
808 analyzed by the trained model increases. These methods are therefore important when a
809 model is intended to be widely deployed or be responsive to an emerging outbreak where
810 policy decisions must be formulated rapidly. Because amortized approximate methods
811 require multiple analyses to realize time savings, researchers need to generate training data
812 sets over a broad parameter and model space so that trained networks can be applied to
813 new and diverse data sets.

814 [R1.1] Our analysis introduces a simple approach to estimate the ancestral state
815 corresponding to the root node or stem node of a phylogeny. More sophisticated supervised
816 learning approaches will be needed to train neural networks to predict the ancestral
817 locations for internal nodes other than the root. The topologies and branch lengths of
818 random phylogenies in the training and test datasets will vary from tree to tree. Our
819 approach relies on the fact that all trees contain a root node, meaning all trees can help
820 predict the root node's state. However, few (if any) trees in the training dataset will contain
821 an arbitrary clade of interest within a test dataset, suggesting to us that naive approaches
822 to train networks to estimate ancestral states for all internal nodes will probably fail. We

823 are unaware of any existing solutions for generalized ancestral state estimation using deep
824 learning, and expect the problem will gather more attention as the field matures.

825 Quantifying uncertainty is crucial to data analysis and decision making, and
826 Bayesian statistics provides a framework for doing so in a rigorous way. [EIC, R1.2]: It is
827 essential to understand how uncertainty estimation with likelihood-free methods compare
828 to likelihood-based methods when confronted with the mismatch of models and real-world
829 data-generating processes. We quantified uncertainty using conformalized quantile
830 regression (CQR; Romano et al. 2019) by training neural networks to predict quantiles and
831 then calibrating those quantiles to produce the expected coverage. We refer to the resulting
832 intervals as calibrated probability intervals (CPI) and demonstrate that they predict well
833 the coverage of true values on a test dataset (Figure 3) and behave in similar ways to
834 Bayesian methods when the model is or is not misspecified (Figures 2 - 9). Despite having
835 the same (correct) coverage as the Bayesian HPI, the interval length was 20-50% wider on
836 average making them a more conservative (less precise) estimation procedure. Though this
837 can likely be improved with more training data for qCNNs, there are more fundamental
838 challenges for uncertainty quantification with quantile regression and conformalization.

839 Methods for estimating more precise intervals is an active vein of research among
840 machine learning researchers and statisticians (Barber et al. 2020; Chung et al. 2021; Sousa
841 et al. 2022; Gibbs et al. 2023). For example, although intervals estimated by the qCNN are
842 conditional on each data point, the calibration of quantiles through conformalization
843 involves estimating marginal calibration terms that shift all quantiles by the same amount.
844 If the error in the quantile coverage is not constant across the prediction range, then a
845 more adaptive procedure should yeild more precise intervals (Sousa et al. 2022; Gibbs et al.
846 2023).

847 We also compared the consistency among CPI estimates at different inner-quantiles
848 to that of HPIS at those same quantiles. We find that independently trained neural
849 networks for each α level can potentially lead to inconsistencies where narrower, nested

850 inner quantiles can have close to or higher coverage than wider quantiles (*e.g.* Figure 9 C).
851 Overall, our results suggest CQR is approximately consistent with likelihood-based
852 methods and similarly sensitive to model misspecification, while there is room for
853 improvement. Methods where all quantiles of interest can be estimated jointly (Chung
854 et al. 2021) may be a fruitful avenue of research for such improvements.

855 Another important challenge of inference with deep learning is the problem of
856 convergence to a location on the loss function surface that approximates the maximum
857 likelihood well. There are a number of basic heuristics that can help such as learning
858 curves but more rigorous methods of ascertaining convergence is the subject of active
859 research (Bürkner et al. 2022; Schmitt et al. 2022).

860 With recent advances in deep learning in epidemiology, evolution, and ecology
861 (Battey et al. 2020; Schrider and Kern 2018; Voznica et al. 2022; Radev et al. 2021;
862 Lambert et al. 2022; Rosenzweig et al. 2022; Suvorov and Schrider 2022) biologists can now
863 explore the behavior of entire classes of stochastic branching models that are biologically
864 interesting but mathematically or statistically prohibitive for use with traditional
865 likelihood-based inference techniques. [EIC] Beyond epidemiology, we anticipate that deep
866 learning approaches will be useful for a wide range of currently intractable phylogenetic
867 modeling problems. Many phylogenetic scenarios – such as the adaptive radiation of anoles
868 (Patton et al. 2021) or the global spread of the grasses (Palazzi et al. 2022) – involve the
869 evolution of discrete traits, continuous traits, speciation, and extinction within an
870 ecological or spatial context across a set of co-evolving species. Deriving fully mechanistic
871 yet tractable phylogenetic model likelihoods for such complex scenarios is difficult, if not
872 impossible. Careful development and applications of likelihood-free modeling methods
873 might bring these phylogenetic scenarios into renewed focus for more detailed study.
874 Although we are cautiously optimistic about the future of deep learning methods for
875 phylogenetics, it will become increasingly important for the field to diagnose the conditions
876 where phylogenetic deep learning underperforms relative to likelihood-based approaches,

877 and to devise general solutions to benefit the field.

878

FUNDING

879 This research was supported in part by an appointment to the Department of Defense
880 (DOD) Research Participation Program administered by the Oak Ridge Institute for
881 Science and Education (ORISE) through an interagency agreement between the U.S.
882 Department of Energy (DOE) and the DOD. ORISE is managed by ORAU under DOE
883 contract number DE-SC0014664. All opinions expressed in this paper are the author's and
884 do not necessarily reflect the policies and views of DOD, DOE, or ORAU/ORISE. MJL
885 was supported by the National Science Foundation (DEB 2040347) and by an internal
886 grant awarded by the Incubator for Transdisciplinary Futures at Washington University.

887

ACKNOWLEDGEMENTS

888 We are grateful to Fábio Mendes, Sarah Swiston, Sean McHugh, Walker Sexton, and
889 Mariana Braga for helpful comments on the research.

890 Data available from the Dryad Digital Repository: <https://doi.org/10.25338/B8SH2J>
891 (Thompson et al. 2023)

893 References

- 894 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
 895 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian
 896 Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal
 897 Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mane, Rajat
 898 Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens,
 899 Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay
 900 Vasudevan, Fernanda Viegas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin
 901 Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-Scale Machine Learning on
 902 Heterogeneous Distributed Systems, March 2016.
- 903 Laith Alzubaidi, Jinglan Zhang, Amjad J. Humaidi, Ayad Al-Dujaili, Ye Duan, Omran
 904 Al-Shamma, J. Santamaría, Mohammed A. Fadhel, Muthana Al-Amidie, and Laith
 905 Farhan. Review of deep learning: Concepts, CNN architectures, challenges, applications,
 906 future directions. *Journal of Big Data*, 8(1):53, 2021. ISSN 2196-1115. doi:
 907 10.1186/s40537-021-00444-8.
- 908 Roy M Anderson and Robert M May. Population biology of infectious diseases: Part i.
 909 *Nature*, 280(5721):361–367, 1979.
- 910 Anastasios N. Angelopoulos, Stephen Bates, Clara Fannjiang, Michael I. Jordan, and
 911 Tijana Zrnic. Prediction-Powered Inference, February 2023.
- 912 Rina Foygel Barber, Emmanuel J. Candès, Aaditya Ramdas, and Ryan J. Tibshirani. The
 913 limits of distribution-free conditional predictive inference, April 2020.
- 914 CJ Battey, Peter L Ralph, and Andrew D Kern. Predicting geographic location from
 915 genetic variation with deep neural networks. *eLife*, 9:e54507, June 2020. ISSN
 916 2050-084X. doi: 10.7554/eLife.54507.

- 917 Jeremy M. Beaulieu and Brian C. O'Meara. Detecting Hidden Diversification Shifts in
918 Models of Trait-Dependent Speciation and Extinction. *Systematic Biology*, 65(4):
919 583–601, July 2016. ISSN 1063-5157, 1076-836X. doi: 10.1093/sysbio/syw022.
- 920 Ottar N. Bjørnstad, Katriona Shea, Martin Krzywinski, and Naomi Altman. The SEIRS
921 model for infectious disease dynamics. *Nature Methods*, 17(6):557–558, June 2020. ISSN
922 1548-7091, 1548-7105. doi: 10.1038/s41592-020-0856-2.
- 923 Folmer Bokma. Artificial neural networks can learn to estimate extinction rates from
924 molecular phylogenies. *Journal of theoretical biology*, 243(3):449–454, 2006.
- 925 Remco Bouckaert, Timothy G. Vaughan, Joëlle Barido-Sottani, Sebastián Duchêne,
926 Mathieu Fourment, Alexandra Gavryushkina, Joseph Heled, Graham Jones, Denise
927 Kühnert, Nicola De Maio, Michael Matschiner, Fábio K. Mendes, Nicola F. Müller,
928 Huw A. Ogilvie, Louis du Plessis, Alex Popinga, Andrew Rambaut, David Rasmussen,
929 Igor Siveroni, Marc A. Suchard, Chieh-Hsi Wu, Dong Xie, Chi Zhang, Tanja Stadler, and
930 Alexei J. Drummond. BEAST 2.5: An advanced software platform for Bayesian
931 evolutionary analysis. *PLOS Computational Biology*, 15(4):e1006650, April 2019. ISSN
932 1553-7358. doi: 10.1371/journal.pcbi.1006650.
- 933 Paul-Christian Bürkner, Maximilian Scholz, and Stefan Radev. Some models are useful,
934 but how do we know which ones? Towards a unified Bayesian model taxonomy,
935 September 2022.
- 936 Sheryl L. Chang, Mahendra Piraveenan, Philippa Pattison, and Mikhail Prokopenko.
937 Game theoretic modelling of infectious disease dynamics and intervention methods: A
938 review. *Journal of Biological Dynamics*, 14(1):57–89, January 2020. ISSN 1751-3758,
939 1751-3766. doi: 10.1080/17513758.2020.1720322.
- 940 F. K. Chollet. Keras: The Python deep learning API. <https://keras.io/>.

- 941 Youngseog Chung, Willie Neiswanger, Ian Char, and Jeff Schneider. Beyond Pinball Loss:
942 Quantile Methods for Calibrated Uncertainty Quantification, December 2021.
- 943 Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based
944 inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062,
945 December 2020. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1912789117.
- 946 Emanuel Masiero da Fonseca, Guarino R. Colli, Fernanda P. Werneck, and Bryan C.
947 Carstens. Phylogeographic model selection using convolutional neural networks,
948 September 2020.
- 949 Jordan Douglas, Fábio K Mendes, Remco Bouckaert, Dong Xie, Cinthy L Jiménez-Silva,
950 Christiaan Swanepoel, Joep de Ligt, Xiaoyun Ren, Matt Storey, James Hadfield, Colin R
951 Simpson, Jemma L Geoghegan, Alexei J Drummond, and David Welch. Phylodynamics
952 reveals the role of human travel and contact tracing in controlling the first wave of
953 COVID-19 in four island nations. *Virus Evolution*, 7(2), September 2021. ISSN
954 2057-1577. doi: 10.1093/ve/veab052.
- 955 Richard G. FitzJohn. Diversitree: Comparative phylogenetic analyses of diversification in
956 R. *Methods in Ecology and Evolution*, 3(6):1084–1092, 2012. ISSN 2041-210X. doi:
957 10.1111/j.2041-210X.2012.00234.x.
- 958 Lex Flagel, Yaniv Brandvain, and Daniel R Schrider. The Unreasonable Effectiveness of
959 Convolutional Neural Networks in Population Genetic Inference. *Molecular Biology and
960 Evolution*, 36(2):220–238, February 2019. ISSN 0737-4038, 1537-1719. doi:
961 10.1093/molbev/msy224.
- 962 Jiansi Gao, Michael R May, Bruce Rannala, and Brian R Moore. New Phylogenetic Models
963 Incorporating Interval-Specific Dispersal Dynamics Improve Inference of Disease Spread.
964 *Molecular Biology and Evolution*, 39(8):msac159, August 2022. ISSN 1537-1719. doi:
965 10.1093/molbev/msac159.

- 966 Jiansi Gao, Michael R. May, Bruce Rannala, and Brian R. Moore. Model misspecification
967 misleads inference of the spatial dynamics of disease outbreaks. *Proceedings of the*
968 *National Academy of Sciences*, 120(11):e2213913120, March 2023. doi:
969 10.1073/pnas.2213913120.
- 970 Isaac Gibbs, John J. Cherian, and Emmanuel J. Candès. Conformal Prediction With
971 Conditional Guarantees, May 2023.
- 972 James Hadfield, Colin Megill, Sidney M Bell, John Huddleston, Barney Potter, Charlton
973 Callender, Pavel Sagulenko, Trevor Bedford, and Richard A Neher. Nextstrain: real-time
974 tracking of pathogen evolution. *Bioinformatics*, 34(23):4121–4123, December 2018. ISSN
975 1367-4803. doi: 10.1093/bioinformatics/bty407. URL
976 <https://doi.org/10.1093/bioinformatics/bty407>.
- 977 Oskar Hagen, Benjamin Flück, Fabian Fopp, Juliano S. Cabral, Florian Hartig, Mikael
978 Pontarp, Thiago F. Rangel, and Loïc Pellissier. Gen3sis: A general engine for
979 eco-evolutionary simulations of the processes that shape Earth’s biodiversity. *PLOS*
980 *Biology*, 19(7):e3001340, July 2021. ISSN 1545-7885. doi: 10.1371/journal.pbio.3001340.
- 981 Benjamin C Haller and Philipp W Messer. SLiM 3: Forward Genetic Simulations Beyond
982 the Wright–Fisher Model. *Molecular Biology and Evolution*, 36(3):632–637, March 2019.
983 ISSN 0737-4038. doi: 10.1093/molbev/msy228.
- 984 Sebastian Höhna, Michael J. Landis, Tracy A. Heath, Bastien Boussau, Nicolas Lartillot,
985 Brian R. Moore, John P. Huelsenbeck, and Fredrik Ronquist. RevBayes: Bayesian
986 Phylogenetic Inference Using Graphical Models and an Interactive Model-Specification
987 Language. *Systematic Biology*, 65(4):726–736, July 2016. ISSN 1063-5157, 1076-836X.
988 doi: 10.1093/sysbio/syw021.
- 989 Eddie C Holmes and Geoff P Garnett. Genes, trees and infections: molecular evidence in
990 epidemiology. *Trends in Ecology & Evolution*, 9(7):256–260, 1994.

- 991 Eddie C Holmes, Sean Nee, Andrew Rambaut, Geoff P Garnett, and Paul H Harvey.
- 992 Revealing the history of infectious disease epidemics through phylogenetic trees.
- 993 *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*,
- 994 349(1327):33–40, 1995.
- 995 Asifullah Khan, Anabia Sohail, Umme Zahoor, and Aqsa Saeed Qureshi. A survey of the
- 996 recent architectures of deep convolutional neural networks. *Artificial Intelligence Review*,
- 997 53(8):5455–5516, December 2020. ISSN 1573-7462. doi: 10.1007/s10462-020-09825-6.
- 998 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization,
- 999 January 2017.
- 1000 Denise Kühnert, Tanja Stadler, Timothy G. Vaughan, and Alexei J. Drummond.
- 1001 Simultaneous reconstruction of evolutionary history and epidemiological dynamics from
- 1002 viral sequences with the birth–death SIR model. *Journal of The Royal Society Interface*,
- 1003 11(94):20131106, May 2014. ISSN 1742-5689, 1742-5662. doi: 10.1098/rsif.2013.1106.
- 1004 Denise Kühnert, Tanja Stadler, Timothy G. Vaughan, and Alexei J. Drummond.
- 1005 Phylodynamics with Migration: A Computational Framework to Quantify Population
- 1006 Structure from Genomic Data. *Molecular Biology and Evolution*, 33(8):2102–2116,
- 1007 August 2016. ISSN 0737-4038. doi: 10.1093/molbev/msw064.
- 1008 Sophia Lambert, Jakub Voznica, and Hélène Morlon. Deep Learning from Phylogenies for
- 1009 Diversification Analyses, September 2022.
- 1010 Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman.
- 1011 Distribution-Free Predictive Inference for Regression. *Journal of the American Statistical*
- 1012 *Association*, 113(523):1094–1111, July 2018. ISSN 0162-1459. doi:
- 1013 10.1080/01621459.2017.1307116.
- 1014 Philippe Lemey, Andrew Rambaut, Alexei J. Drummond, and Marc A. Suchard. Bayesian

- 1015 Phylogeography Finds Its Roots. *PLoS Computational Biology*, 5(9):e1000520,
1016 September 2009. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1000520.
- 1017 Philippe Lemey, Nick Ruktanonchai, Samuel L. Hong, Vittoria Colizza, Chiara Poletto,
1018 Frederik Van den Broeck, Mandev S. Gill, Xiang Ji, Anthony Levasseur, Bas B.
1019 Oude Munnink, Marion Koopmans, Adam Sadilek, Shengjie Lai, Andrew J. Tatem, Guy
1020 Baele, Marc A. Suchard, and Simon Dellicour. Untangling introductions and persistence
1021 in COVID-19 resurgence in Europe. *Nature*, June 2021. ISSN 0028-0836, 1476-4687. doi:
1022 10.1038/s41586-021-03754-2.
- 1023 Frédéric Lemoine and Olivier Gascuel. Gotree/Goalign: Toolkit and Go API to facilitate
1024 the development of phylogenetic workflows. *NAR Genomics and Bioinformatics*, 3(3):
1025 lqab075, September 2021. ISSN 2631-9268. doi: 10.1093/nargab/lqab075.
- 1026 Wayne P. Maddison, Peter E. Midford, and Sarah P. Otto. Estimating a Binary
1027 Character's Effect on Speciation and Extinction. *Systematic Biology*, 56(5):701–710,
1028 October 2007. ISSN 1076-836X, 1063-5157. doi: 10.1080/10635150701607033.
- 1029 Mike Meredith and John Kruschke. Bayesian Estimation Supersedes the t-Test. page 13.
- 1030 Bui Quang Minh, Heiko A Schmidt, Olga Chernomor, Dominik Schrempf, Michael D
1031 Woodhams, Arndt von Haeseler, and Robert Lanfear. IQ-TREE 2: New Models and
1032 Efficient Methods for Phylogenetic Inference in the Genomic Era. *Molecular Biology and*
1033 *Evolution*, 37(5):1530–1534, May 2020. ISSN 0737-4038. doi: 10.1093/molbev/msaa015.
- 1034 Niema Moshiri, Manon Ragonnet-Cronin, Joel O Wertheim, and Siavash Mirarab.
1035 FAVITES: Simultaneous simulation of transmission networks, phylogenetic trees and
1036 sequences. *Bioinformatics*, 35(11):1852–1861, June 2019. ISSN 1367-4803, 1460-2059.
1037 doi: 10.1093/bioinformatics/bty921.
- 1038 Sarah A. Nadeau, Timothy G. Vaughan, Jérémie Scire, Jana S. Huisman, and Tanja
1039 Stadler. The origin and early spread of SARS-CoV-2 in Europe. *Proceedings of the*

- 1040 National Academy of Sciences, 118(9):e2012008118, March 2021. ISSN 0027-8424,
1041 1091-6490. doi: 10.1073/pnas.2012008118.
- 1042 Luca Nesterenko, Bastien Boussau, and Laurent Jacob. Phyloformer: Towards fast and
1043 accurate phylogeny estimation with self-attention networks, June 2022.
- 1044 Eamon B O'Dea and John M Drake. A semi-parametric, state-space compartmental model
1045 with time-dependent parameters for forecasting COVID-19 cases, hospitalizations, and
1046 deaths. page 32, 2021.
- 1047 Isaac Overcast, Megan Ruffley, James Rosindell, Luke Harmon, Paulo AV Borges, Brent C
1048 Emerson, Rampal S Etienne, Rosemary Gillespie, Henrik Krehenwinkel, D Luke Mahler,
1049 et al. A unified model of species abundance, genetic diversity, and functional diversity
1050 reveals the mechanisms structuring ecological communities. *Molecular Ecology*
1051 *Resources*, 21(8):2782–2800, 2021.
- 1052 Luis Palazzi, Oriane Hidalgo, Viviana D Barreda, Félix Forest, and Sebastian Höhna.
1053 The rise of grasslands is linked to atmospheric co₂ decline in the late palaeogene. *Nature*
1054 *Communications*, 13:293, 2022.
- 1055 Austin H Patton, Luke J Harmon, María del Rosario Castañeda, Hannah K Frank, Colin M
1056 Donihue, Anthony Herrel, and Jonathan B Losos. When adaptive radiations collide:
1057 Different evolutionary trajectories between and within island and mainland lizard clades.
1058 *Proceedings of the National Academy of Sciences*, 118(42):e2024451118, 2021.
- 1059 Jonathan E. Pekar, Andrew Magee, Edyth Parker, Niema Moshiri, Katherine Izhikevich,
1060 Jennifer L. Havens, Karthik Gangavarapu, Lorena Mariana Malpica Serrano, Alexander
1061 Crits-Christoph, Nathaniel L. Matteson, Mark Zeller, Joshua I. Levy, Jade C. Wang,
1062 Scott Hughes, Jungmin Lee, Heedo Park, Man-Seong Park, Katherine Zi Yan Ching,
1063 Raymond Tzer Pin Lin, Mohd Noor Mat Isa, Yusuf Muhammad Noor, Tetyana I.
1064 Vasylyeva, Robert F. Garry, Edward C. Holmes, Andrew Rambaut, Marc A. Suchard,

- 1065 Kristian G. Andersen, Michael Worobey, and Joel O. Wertheim. The molecular
1066 epidemiology of multiple zoonotic origins of SARS-CoV-2. *Science*, 0(0):eabp8337, July
1067 2022. doi: 10.1126/science.abp8337.
- 1068 José M. Ponciano and Marcos A. Capistrán. First Principles Modeling of Nonlinear
1069 Incidence Rates in Seasonal Epidemics. *PLOS Computational Biology*, 7(2):e1001079,
1070 February 2011. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1001079.
- 1071 O. G. Pybus, M. A. Suchard, P. Lemey, F. J. Bernardin, A. Rambaut, F. W. Crawford,
1072 R. R. Gray, N. Arinaminpathy, S. L. Stramer, M. P. Busch, and E. L. Delwart. Unifying
1073 the spatial epidemiology and molecular evolution of emerging epidemics. *Proceedings of
1074 the National Academy of Sciences*, 109(37):15066–15071, September 2012. ISSN
1075 0027-8424, 1091-6490. doi: 10.1073/pnas.1206598109.
- 1076 Stefan T. Radev, Frederik Graw, Simiao Chen, Nico T. Mutters, Vanessa M. Eichel, Till
1077 Bärnighausen, and Ullrich Köthe. OutbreakFlow: Model-based Bayesian inference of
1078 disease outbreak dynamics with invertible neural networks and its application to the
1079 COVID-19 pandemics in Germany. *PLOS Computational Biology*, 17(10):e1009472,
1080 October 2021. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1009472.
- 1081 A. Rambaut and N. C. Grassly. Seq-Gen: an application for the Monte Carlo simulation of
1082 DNA sequence evolution along phylogenetic trees. *Computer Applications in the
1083 Biosciences*, 13:235–238, 1997.
- 1084 Andrew Rambaut, Oliver G Pybus, Martha I Nelson, Cecile Viboud, Jeffery K
1085 Taubenberger, and Edward C Holmes. The genomic and epidemiological dynamics of
1086 human influenza a virus. *Nature*, 453(7195):615–619, 2008.
- 1087 Liam J. Revell. Phytools: An R package for phylogenetic comparative biology (and other
1088 things). *Methods in Ecology and Evolution*, 3(2):217–223, 2012. ISSN 2041-210X. doi:
1089 10.1111/j.2041-210X.2011.00169.x.

- 1090 Francisco Richter, Bart Haegeman, Rampal S. Etienne, and Ernst C. Wit. Introducing a
1091 general class of species diversification models for phylogenetic trees. *Statistica*
1092 *Neerlandica*, 74(3):261–274, 2020. ISSN 1467-9574. doi: 10.1111/stan.12205.
- 1093 Yaniv Romano, Evan Patterson, and Emmanuel Candès. Conformalized Quantile
1094 Regression. In *Advances in Neural Information Processing Systems*, volume 32. Curran
1095 Associates, Inc., 2019.
- 1096 Benjamin K. Rosenzweig, Matthew W. Hahn, and Andrew Kern. Accurate Detection of
1097 Incomplete Lineage Sorting via Supervised Machine Learning, November 2022.
- 1098 Marvin Schmitt, Paul-Christian Bürkner, Ullrich Köthe, and Stefan T. Radev. Detecting
1099 Model Misspecification in Amortized Bayesian Inference with Neural Networks, May
1100 2022.
- 1101 Daniel R. Schrider and Andrew D. Kern. Supervised Machine Learning for Population
1102 Genetics: A New Paradigm. *Trends in Genetics*, 34(4):301–312, April 2018. ISSN
1103 01689525. doi: 10.1016/j.tig.2017.12.005.
- 1104 Jérémie Scire, Joëlle Barido-Sottani, Denise Kühnert, Timothy G. Vaughan, and Tanja
1105 Stadler. Improved multi-type birth-death phylodynamic inference in BEAST 2. Preprint,
1106 Evolutionary Biology, January 2020.
- 1107 Claudia Solis-Lemus, Shengwen Yang, and Leonardo Zepeda-Nunez. Accurate Phylogenetic
1108 Inference with a Symmetry-preserving Neural Network Model, January 2022.
- 1109 Martim Sousa, Ana Maria Tomé, and José Moreira. Improved conformalized quantile
1110 regression, November 2022.
- 1111 Tanja Stadler. Sampling-through-time in birth–death trees. *Journal of Theoretical Biology*,
1112 267(3):396–404, December 2010. ISSN 00225193. doi: 10.1016/j.jtbi.2010.09.010.

- 1113 Tanja Stadler, Roger Kouyos, Viktor von Wyl, Sabine Yerly, Jürg Böni, Philippe Bürgisser,
1114 Thomas Klimkait, Beda Joos, Philip Rieder, Dong Xie, Huldrych F. Günthard, Alexei J.
1115 Drummond, Sebastian Bonhoeffer, and the Swiss HIV Cohort Study. Estimating the
1116 Basic Reproductive Number from Viral Sequence Data. *Molecular Biology and Evolution*,
1117 29(1):347–357, January 2012. ISSN 1537-1719, 0737-4038. doi: 10.1093/molbev/msr217.
- 1118 Ingo Steinwart and Andreas Christmann. Estimating conditional quantiles with the help of
1119 the pinball loss. *Bernoulli*, 17(1):211–225, February 2011. ISSN 1350-7265. doi:
1120 10.3150/10-BEJ267.
- 1121 Anton Suvorov and Daniel R. Schrider. Reliable estimation of tree branch lengths using
1122 deep neural networks. *bioRxiv*, 2022. doi: 10.1101/2022.11.07.515518. URL
1123 <https://www.biorxiv.org/content/early/2023/02/21/2022.11.07.515518>.
- 1124 Anton Suvorov, Joshua Hochuli, and Daniel R Schrider. Accurate Inference of Tree
1125 Topologies from Multiple Sequence Alignments Using Deep Learning. *Systematic Biology*,
1126 69(2):221–233, March 2020. ISSN 1063-5157, 1076-836X. doi: 10.1093/sysbio/syz060.
- 1127 Ammon Thompson, Benjamin Liebeskind, Erik J. Scully, and Michael J. Landis. Deep
1128 learning phylogeography. *Dryad*, 2023. doi: 10.25338/B8SH2J.
- 1129 Timothy G. Vaughan and Alexei J. Drummond. A Stochastic Simulator of Birth–Death
1130 Master Equations with Application to Phylodynamics. *Molecular Biology and Evolution*,
1131 30(6):1480–1493, June 2013. ISSN 0737-4038. doi: 10.1093/molbev/mst057.
- 1132 Timothy G. Vaughan, Denise Kühnert, Alex Popinga, David Welch, and Alexei J.
1133 Drummond. Efficient Bayesian inference under the structured coalescent.
1134 *Bioinformatics*, 30(16):2272–2279, August 2014. ISSN 1367-4803, 1460-2059. doi:
1135 10.1093/bioinformatics/btu201.
- 1136 Erik M. Volz and Igor Siveroni. Bayesian phylodynamic inference with complex models.

- ₁₁₃₇ *PLOS Computational Biology*, 14(11):e1006546, November 2018. ISSN 1553-7358. doi:
₁₁₃₈ 10.1371/journal.pcbi.1006546.
- ₁₁₃₉ Erik M. Volz, Katia Koelle, and Trevor Bedford. Viral Phylodynamics. *PLOS
1140 Computational Biology*, 9(3):e1002947, March 2013. ISSN 1553-7358. doi:
₁₁₄₁ 10.1371/journal.pcbi.1002947. URL <https:////journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002947>.
₁₁₄₂ Publisher: Public Library of Science.
- ₁₁₄₃ Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. Conformal Prediction: General
₁₁₄₄ Case and Regression. In Vladimir Vovk, Alexander Gammerman, and Glenn Shafer,
₁₁₄₅ editors, *Algorithmic Learning in a Random World*, pages 19–69. Springer International
₁₁₄₆ Publishing, Cham, 2022. ISBN 978-3-031-06649-8. doi: 10.1007/978-3-031-06649-8_2.
- ₁₁₄₇ J Voznica, A Zhukova, V Boskova, E Saulnier, F Lemoine, M Moslonka-Lefebvre, and
₁₁₄₈ O Gascuel. Deep learning from phylogenies to uncover the epidemiological dynamics of
₁₁₄₉ outbreaks. preprint, Bioinformatics, March 2021. URL
₁₁₅₀ <http://biorxiv.org/lookup/doi/10.1101/2021.03.11.435006>.
- ₁₁₅₁ J. Voznica, A. Zhukova, V. Boskova, E. Saulnier, F. Lemoine, M. Moslonka-Lefebvre, and
₁₁₅₂ O. Gascuel. Deep learning from phylogenies to uncover the epidemiological dynamics of
₁₁₅₃ outbreaks. *Nature Communications*, 13(1):3896, July 2022. ISSN 2041-1723. doi:
₁₁₅₄ 10.1038/s41467-022-31511-0.
- ₁₁₅₅ Nicole L. Washington, Karthik Gangavarapu, Mark Zeller, Alexandre Bolze, Elizabeth T.
₁₁₅₆ Cirulli, Kelly M. Schiabor Barrett, Brendan B. Larsen, Catelyn Anderson, Simon White,
₁₁₅₇ Tyler Cassens, Sharoni Jacobs, Geraint Levan, Jason Nguyen, Jimmy M. Ramirez,
₁₁₅₈ Charlotte Rivera-Garcia, Efren Sandoval, Xueqing Wang, David Wong, Emily Spencer,
₁₁₅₉ Refugio Robles-Sikisaka, Ezra Kurzban, Laura D. Hughes, Xianding Deng, Candace
₁₁₆₀ Wang, Venice Servellita, Holly Valentine, Peter De Hoff, Phoebe Seaver, Shashank Sathe,
₁₁₆₁ Wang, Venice Servellita, Holly Valentine, Peter De Hoff, Phoebe Seaver, Shashank Sathe,

- 1162 Kimberly Gietzen, Brad Sickler, Jay Antico, Kelly Hoon, Jingtao Liu, Aaron Harding,
1163 Omid Bakhtar, Tracy Basler, Brett Austin, Duncan MacCannell, Magnus Isaksson,
1164 Phillip G. Febbo, David Becker, Marc Laurent, Eric McDonald, Gene W. Yeo, Rob
1165 Knight, Louise C. Laurent, Eileen de Feo, Michael Worobey, Charles Y. Chiu, Marc A.
1166 Suchard, James T. Lu, William Lee, and Kristian G. Andersen. Emergence and rapid
1167 transmission of SARS-CoV-2 B.1.1.7 in the United States. *Cell*, 184(10):2587–2594.e7,
1168 May 2021. ISSN 00928674. doi: 10.1016/j.cell.2021.03.052.
- 1169 Michael Worobey, Thomas D Watts, Richard A McKay, Marc A Suchard, Timothy
1170 Granade, Dirk E Teuwen, Beryl A Kobrin, Walid Heneine, Philippe Lemey, and
1171 Harold W Jaffe. 1970s and ‘patient 0’hiv-1 genomes illuminate early hiv/aids history in
1172 north america. *Nature*, 539(7627):98–101, 2016.
- 1173 Michael Worobey, Jonathan Pekar, Brendan B. Larsen, Martha I. Nelson, Verity Hill,
1174 Jeffrey B. Joy, Andrew Rambaut, Marc A. Suchard, Joel O. Wertheim, and Philippe
1175 Lemey. The emergence of SARS-CoV-2 in Europe and North America. *Science*, 370
1176 (6516):564–570, October 2020. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.abc8169.

SUPPLEMENTAL TABLES

Table S1: BEST comparisons between CNN and Bayesian absolute percent errors (APEs) for model parameters across all experiments.

95% HPD intervals of average relative error from BEST analysis			
True inference model (Reference for misspecification experiments)	Median CNN APE	Median Like.-based APE	median(CNN APE - Like.-based APE)
R_0	2.4, 3.5	2.1, 3.1	0.1, 1.2
δ	7.0, 10.5	5.7, 8.9	0.2, 3.0
m	9.5, 14.1	8.4, 12.1	0.4, 3.2
<hr/>			
Misspecified R_0 experiment	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	11.8, 17.8	11.0, 16.9	-0.1, 1.6
δ	0.8, 7.6	-0.6, 5.3	1.3, 5.8
m	8.2, 17.9	6.5, 15.9	1.3, 4.7
<hr/>			
Misspecified sample rate experiment	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	-0.3, 1.7	0.03, 1.7	0.1, 1.3
δ	12.0, 21.2	12.6, 21.4	0.1, 4.0
m	3.3, 12.0	5.6, 14.4	-1.2, 2.7
<hr/>			
Misspecified migration rate experiment	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	-0.9, 0.8	-0.6, 1.0	-0.5, 0.8
δ	-2.3, 3.3	0.1, 5.8	-1.4, 2.3
m	4.0, 15.2	5.0, 16.2	-1.3, 2.6
<hr/>			
Misspecified number of locations experiment	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	-0.3, 1.5	-0.7, 0.8	0.5, 1.9
δ	-0.3, 4.9	-0.5, 4.2	0.4, 3.5
m	3.4, 11.1	5.8, 13.5	-0.9, 1.6
<hr/>			
Phylogenetic error experiment	Median CNN APE - median CNN Reference APE	Median Like-based APE - median Reference Like-based APE	median(CNN APE - Like-based APE)
R_0	0.7, 3.0	1.7, 4.4	-1.4, 0.1
δ	2.3, 9.6	1.5, 7.2	1.4, 5.3
m	-1.2, 6.0	-1.8, 5.4	-1.7, 2.4

Table S2: Comparison 95% CPI and HPI for all experiments.

Coverage, width, and overlap of 95% Intervals	\mathbf{R}_0				δ				m			
	CNN CPI	Bayes HPI	Mean CPI width / HPI width	Mean Jaccard index	CNN CPI	Baye s HPI	Mean CPI width / HPI width	Mean Jaccard index	CNN CPI	Bayes HPI	Mean CPI width / HPI width	Mean Jaccard index
True model	0.95	0.96	1.4	0.67	0.96	0.94	1.4	0.66	0.94	0.95	1.2	0.75
Misspecified R_0	0.44	0.29	1.5	0.63	0.9	0.90	1.5	0.63	0.74	0.67	1.2	0.76
Misspecified δ	0.95	0.96	1.4	0.67	0.71	0.55	1.3	0.69	0.72	0.75	1.2	0.75
Misspecified m	0.93	0.94	1.5	0.63	0.94	0.96	1.5	0.65	0.73	0.69	1.3	0.73
Misspecified. Number of locations	0.93	0.92	1.4	0.65	0.96	0.96	1.4	0.68	0.82	0.80	1.2	0.76
Phylogenetic error	0.79	0.60	1.4	0.59	0.71	0.81	1.3	0.59	0.87	0.83	1.3	0.71

SUPPLEMENTAL FIGURES

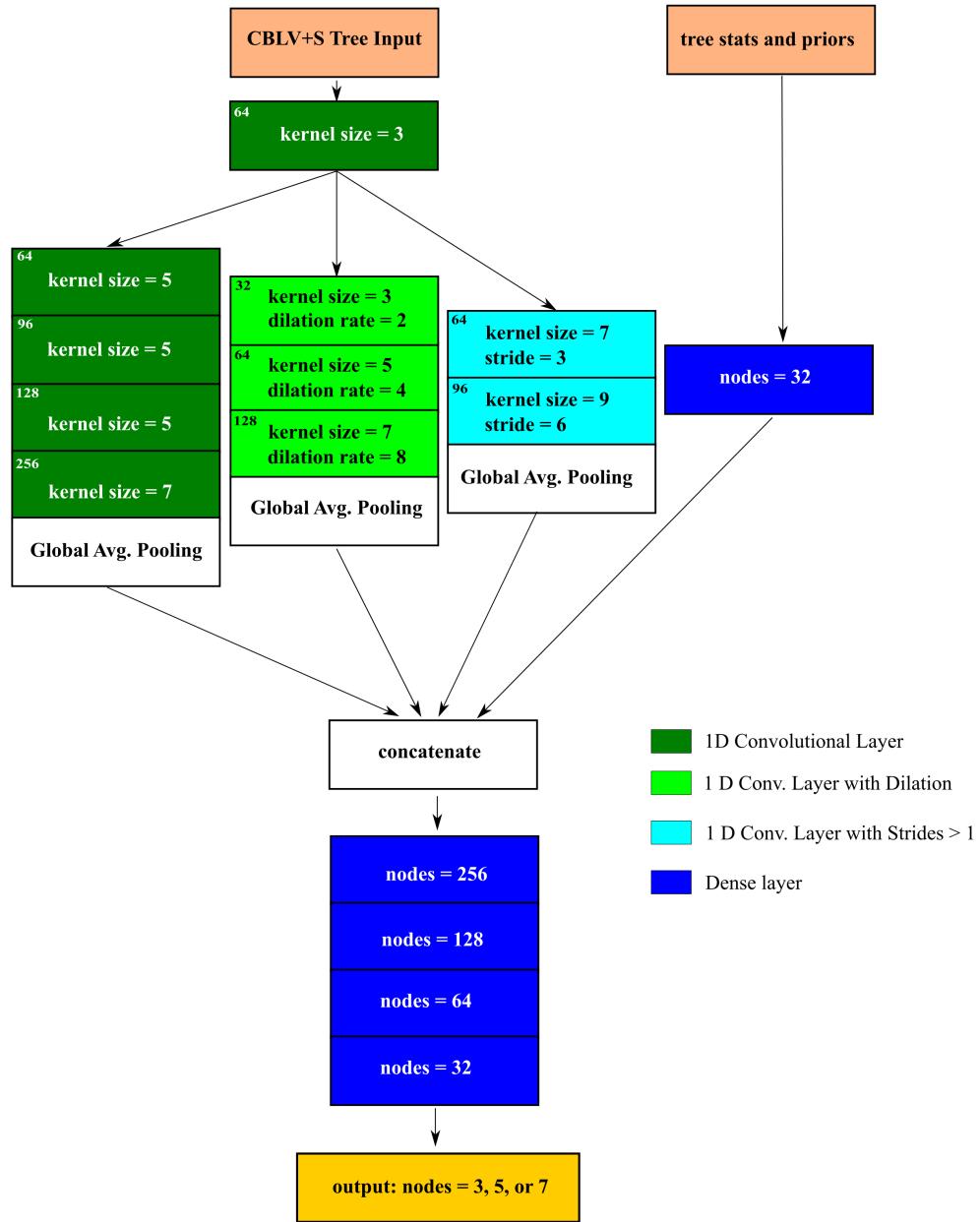


Figure S1: Diagram of deep neural network trained to make 2 kinds of predictions (rates and origin location) under two models (LIBDS and LDBDS).

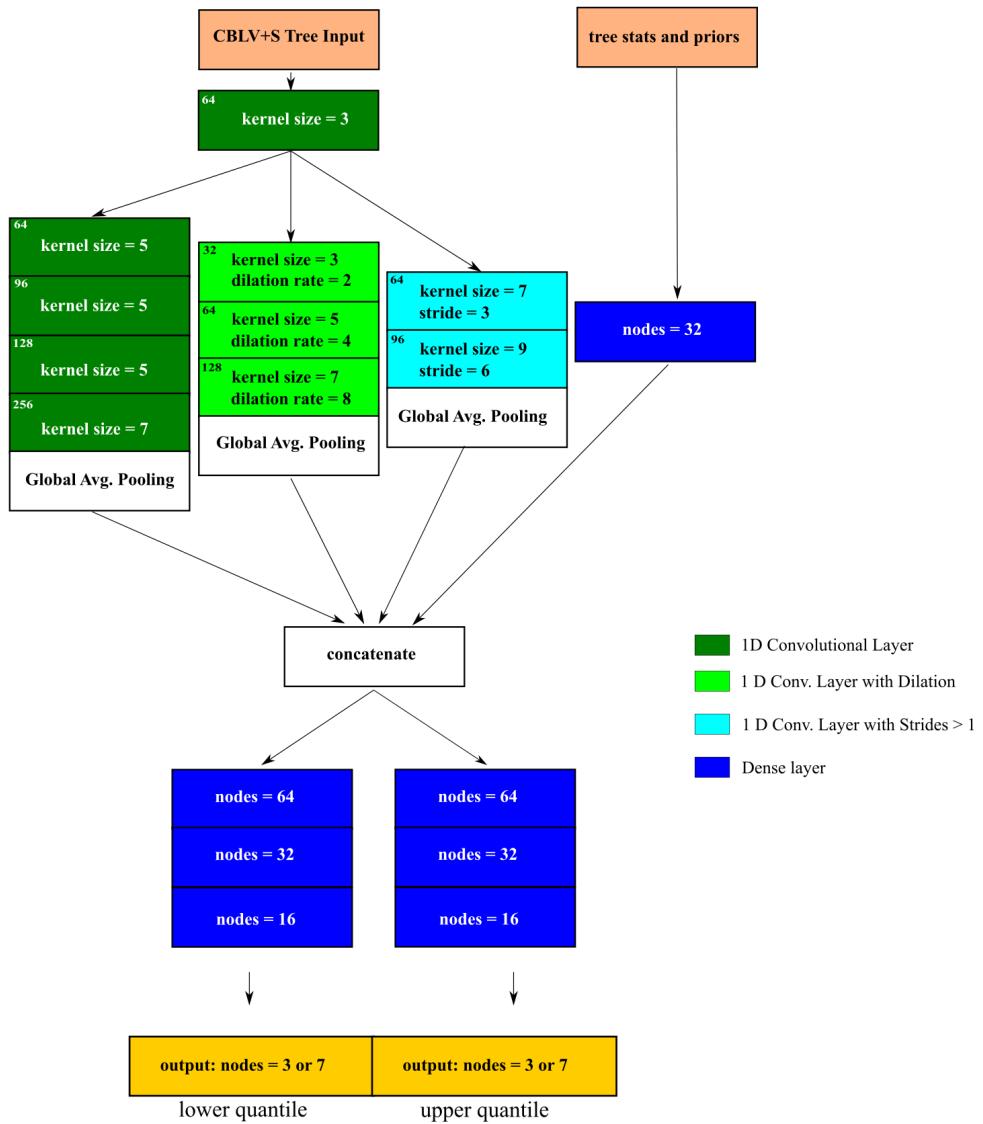


Figure S2: Diagram of deep neural network trained to predict the upper and lower quantiles for a specified α level under two models (LIBDS and LDBDS).

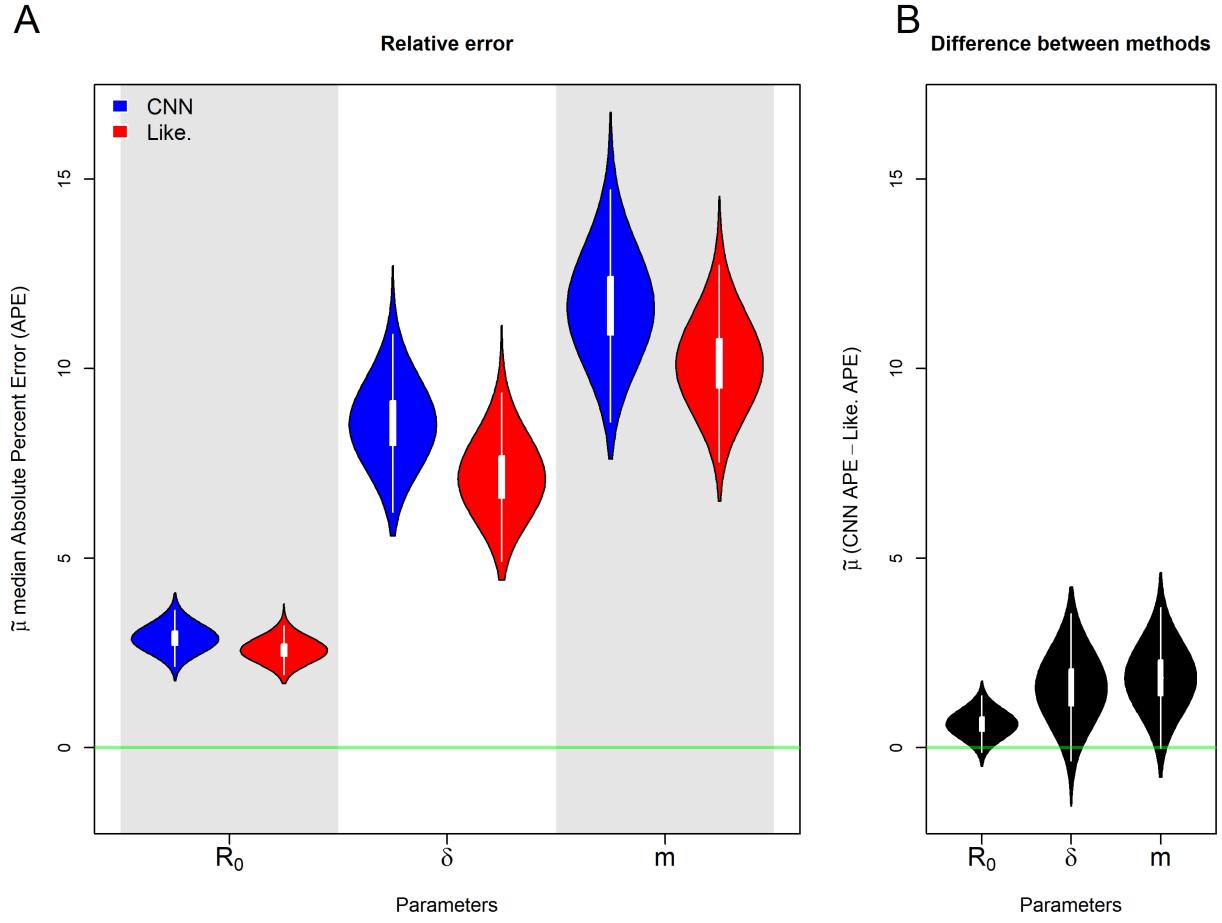


Figure S3: Posterior distributions of the population median, $\tilde{\mu}$, APE estimates of the rate parameters R_0 , δ , and m under the true model. A) shows posterior distribution of the median APE for each of the 3 rate parameters estimated by the CNN (blue) and the likelihood-based method (red). The green line indicates no error. B) shows the posterior distribution for the median difference between the CNN estimate's APE and the likelihood-based estimate's APE. The green line indicates the median APE difference is zero.

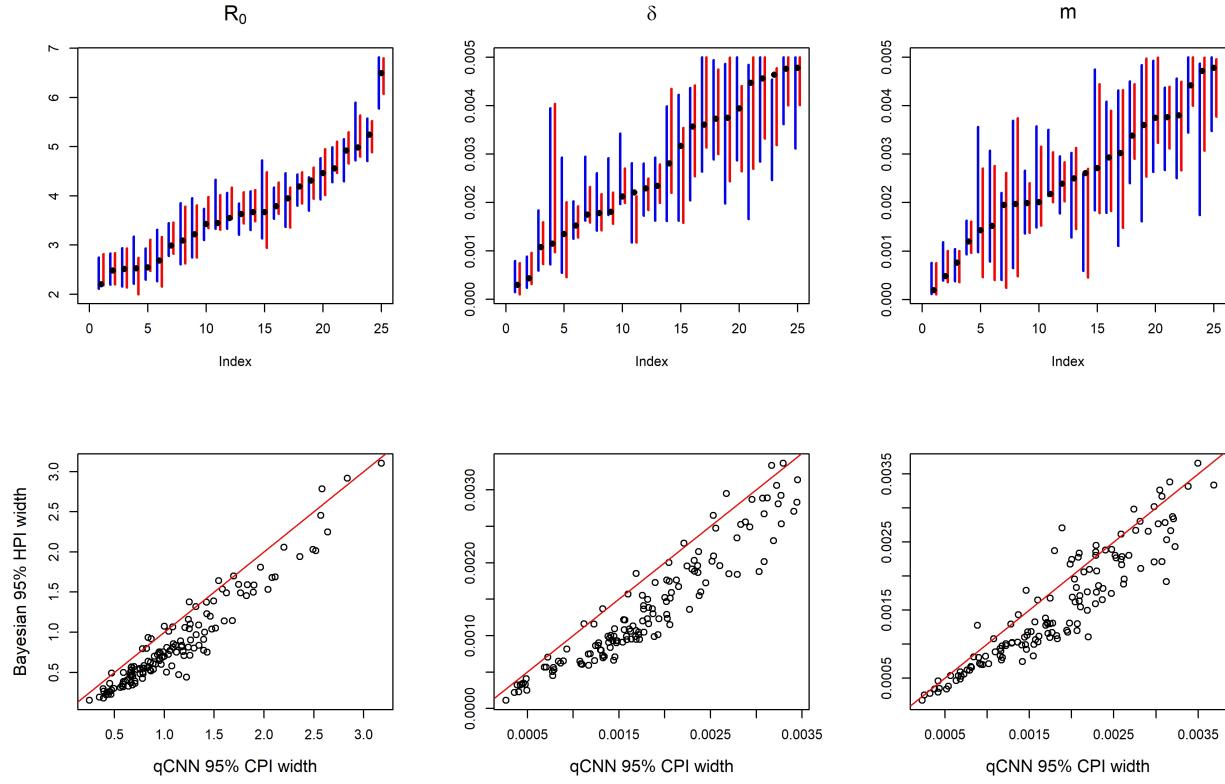


Figure S4: Comparison of interval overlap and relative widths of qCNN and Bayesian methods of uncertainty quantification under the true simulating model. Top row: 95% calibrated probability intervals (CPI) from CNN conformalized quantile regression (blue) and 95% highest posterior intervals (HPI) from Bayesian phylogenetic analysis (red) from a random subset of the data for visualization purposes. Bottom row: scatterplots of the lengths of CPI and HPI intervals of all experiment data. The red diagonal $y = x$ line is for reference.

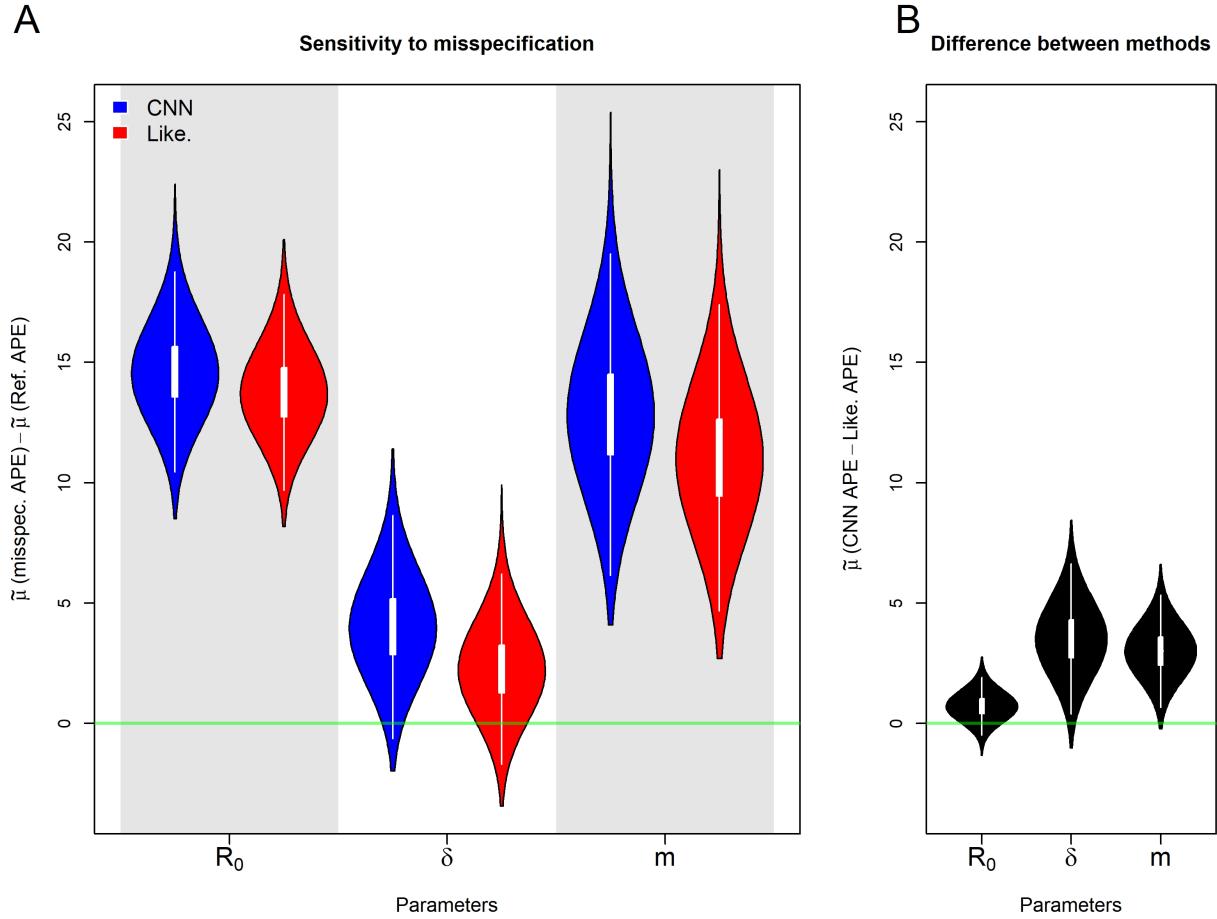


Figure S5: Posterior distributions of median, $\tilde{\mu}$, APE for the misspecified R_0 experiment. A) shows posterior distribution of the difference between the median error under the misspecified model and the the median error under the true, reference model. B) shows the posterior distribution for the population median difference between the CNN estimate's APE and the likelihood-based estimate's APE.

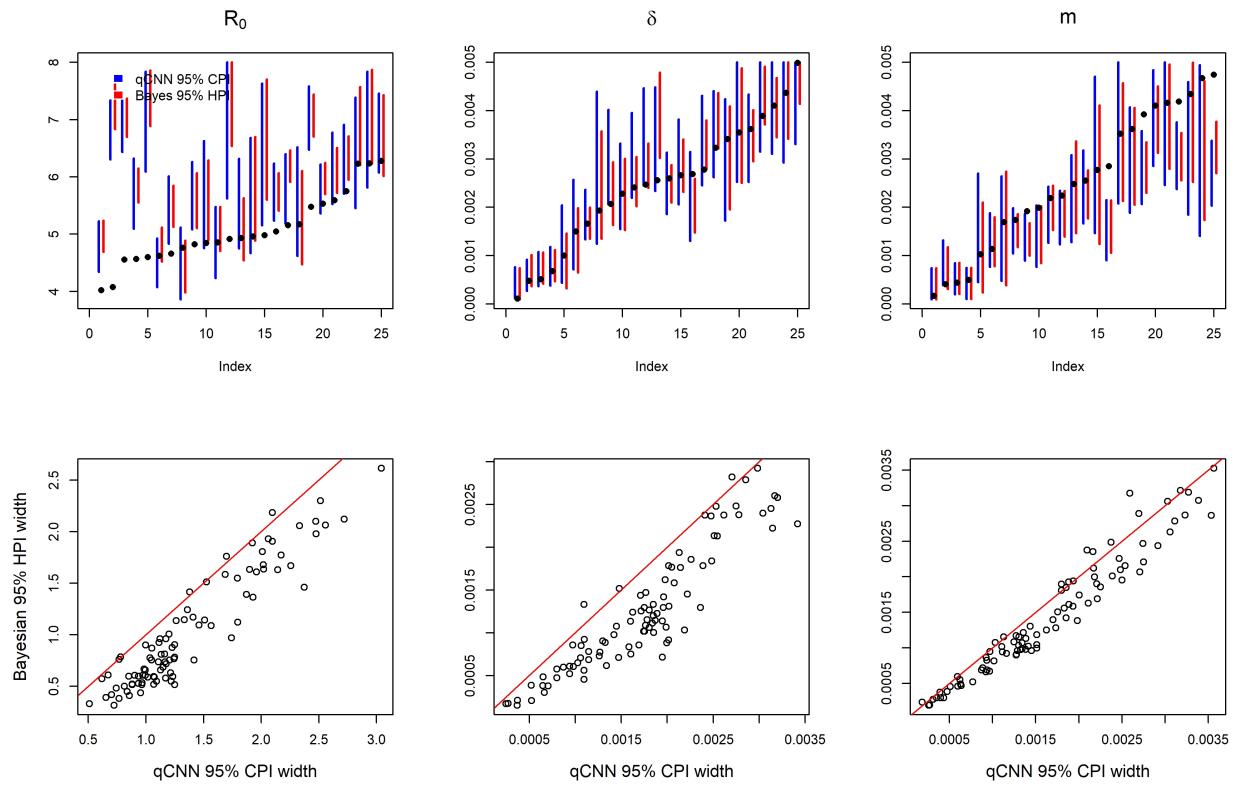


Figure S6: Comparison of CPI and HPI intervals for misspecified R_0 experiment. See SI Figure S4 for general details about plot.

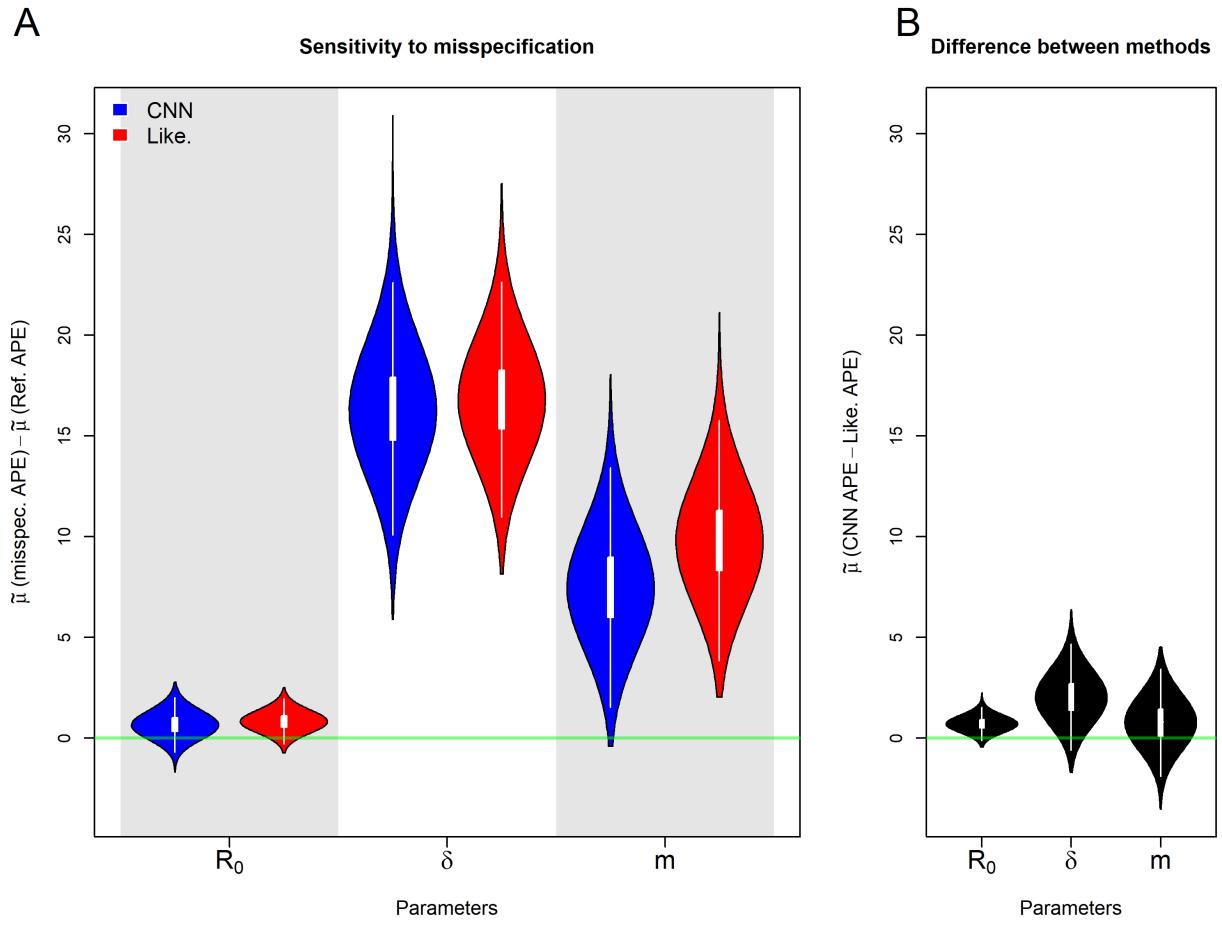


Figure S7: Posterior distributions of median, $\tilde{\mu}$, APE for the misspecified sampling rate, δ , experiment. Details are the same as in S5

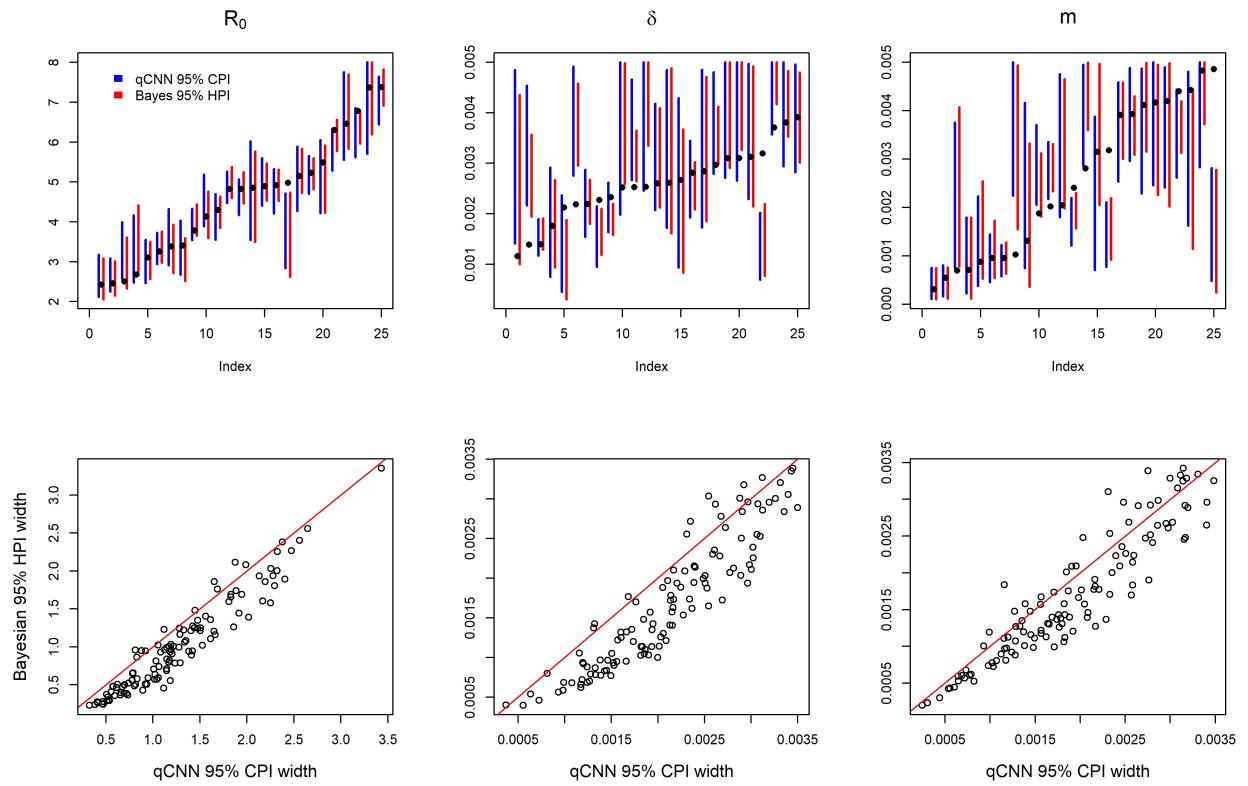


Figure S8: Comparison of CPI and HPI intervals for misspecified δ experiment. See SI Figure S4 for general details about plot.

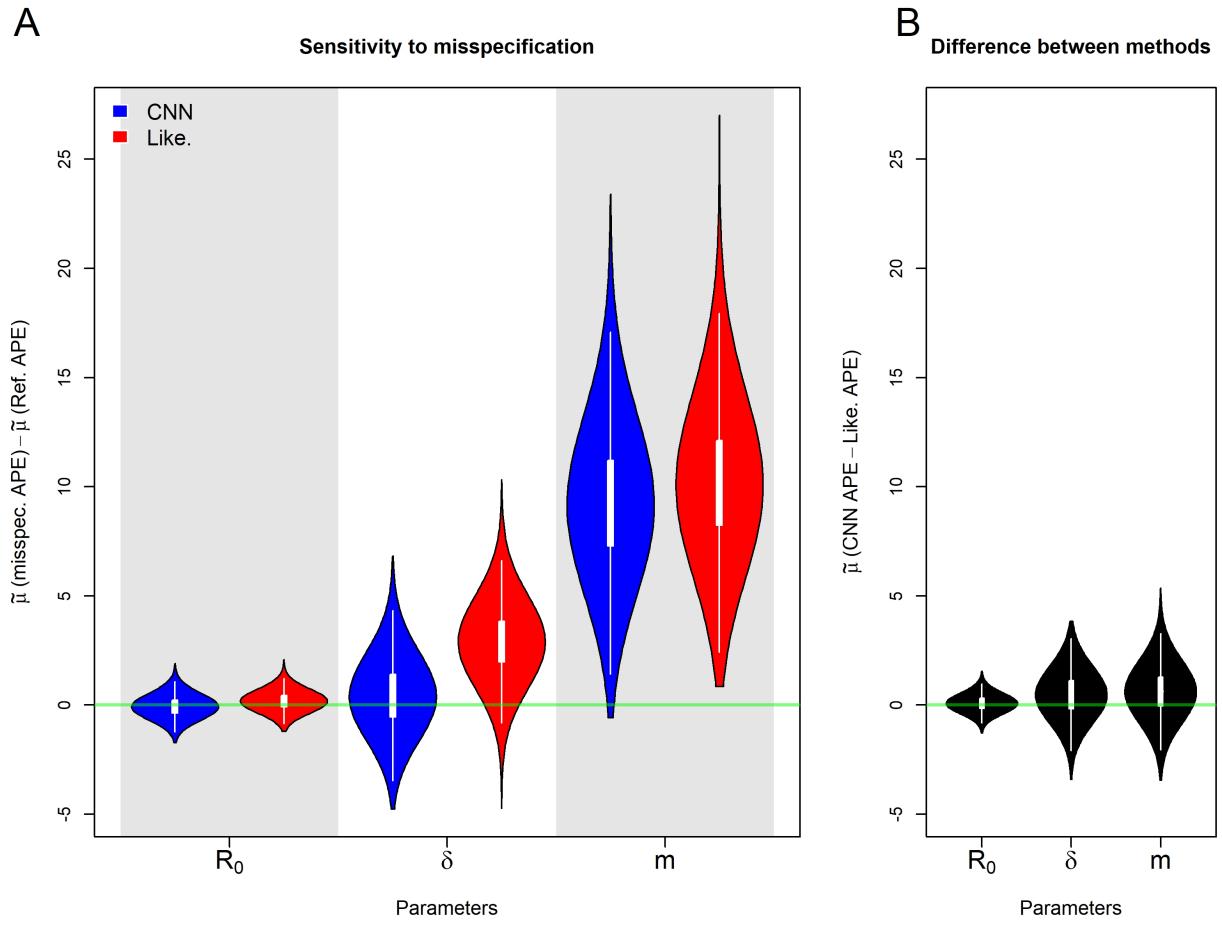


Figure S9: Posterior distributions of median, $\tilde{\mu}$, APE for the misspecified migration rate, m , experiment. Details are the same as in S5

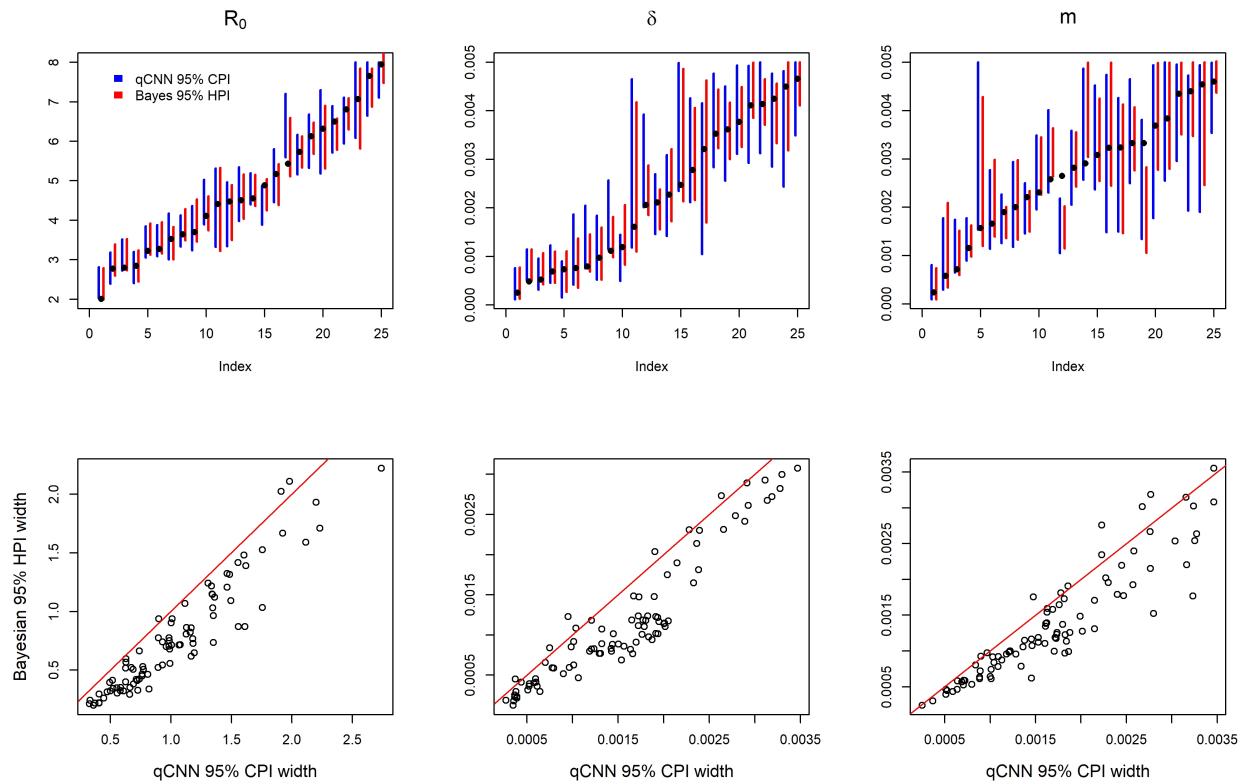


Figure S10: Comparison of CPI and HPI intervals for misspecified migration rate experiment. See SI Figure S4 for general details about plot.

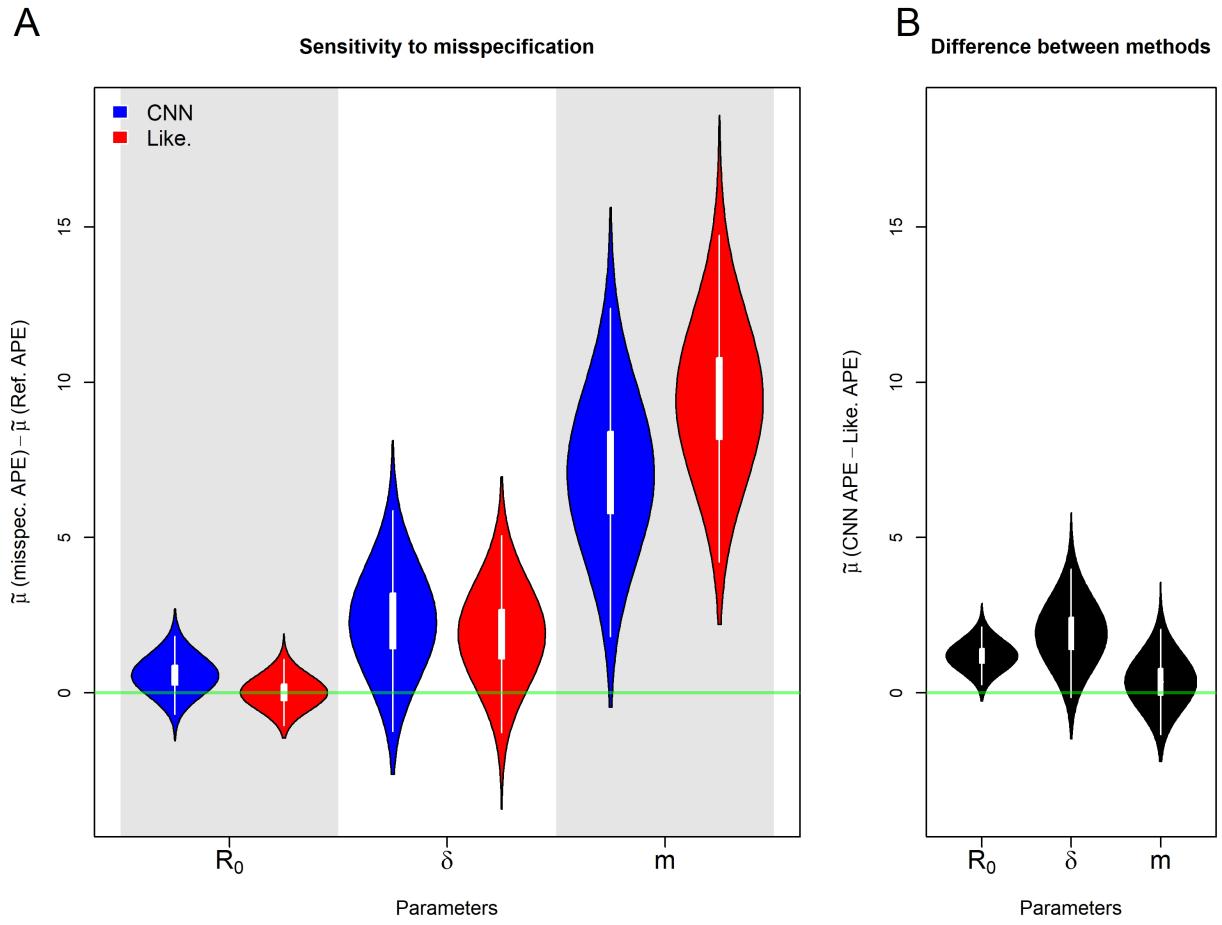


Figure S11: Posterior distributions of the median APE when the model is misspecified for the number of locations. Details are the same as in S5

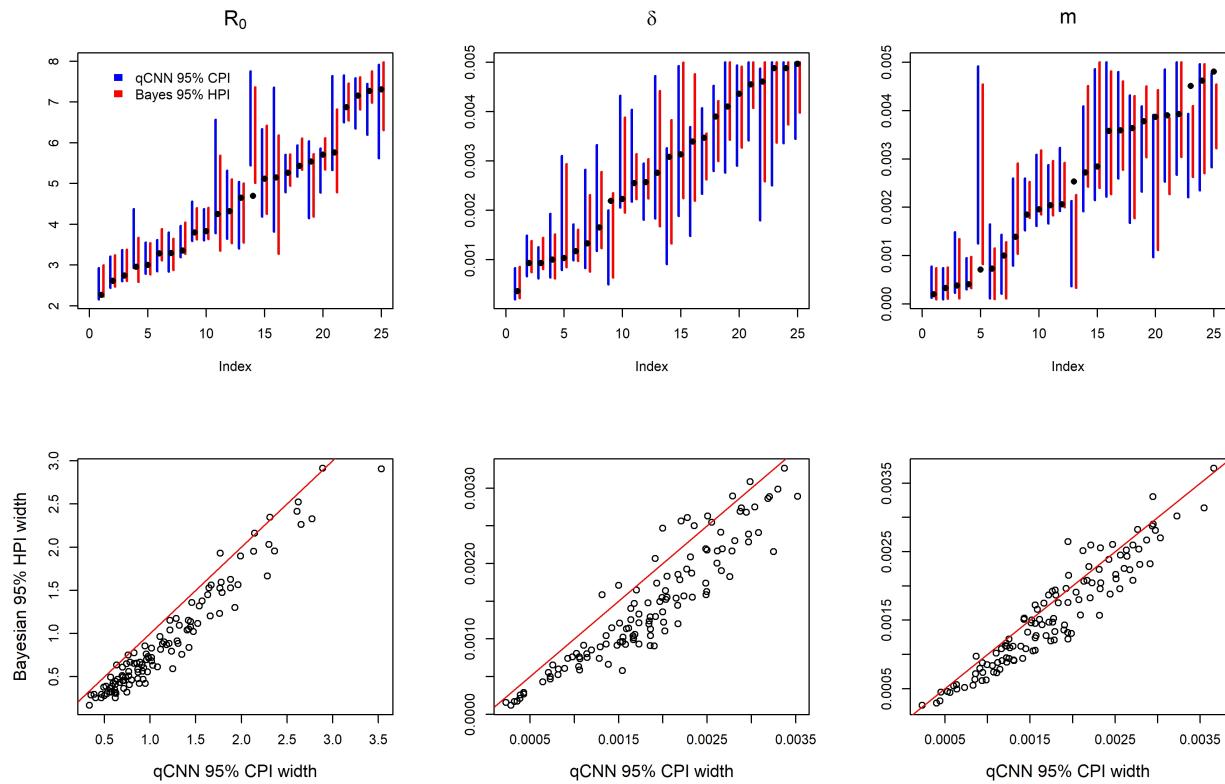


Figure S12: Comparison of CPI and HPI intervals for misspecified number of locations experiment. See SI Figure S4 for general details about plot.

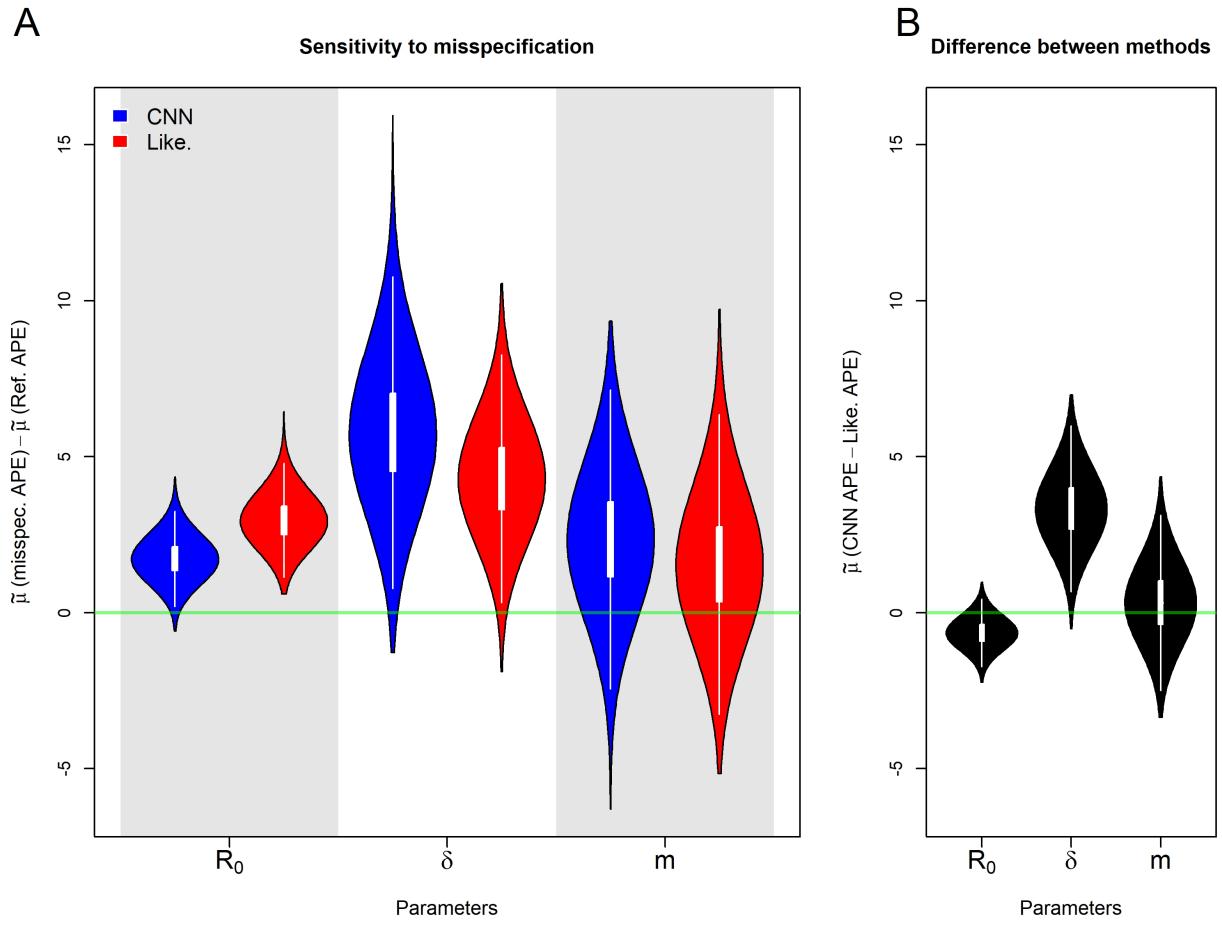


Figure S13: Posterior distributions of the median APE when the phylogenetic tree is incorrect. Details are the same as in S5

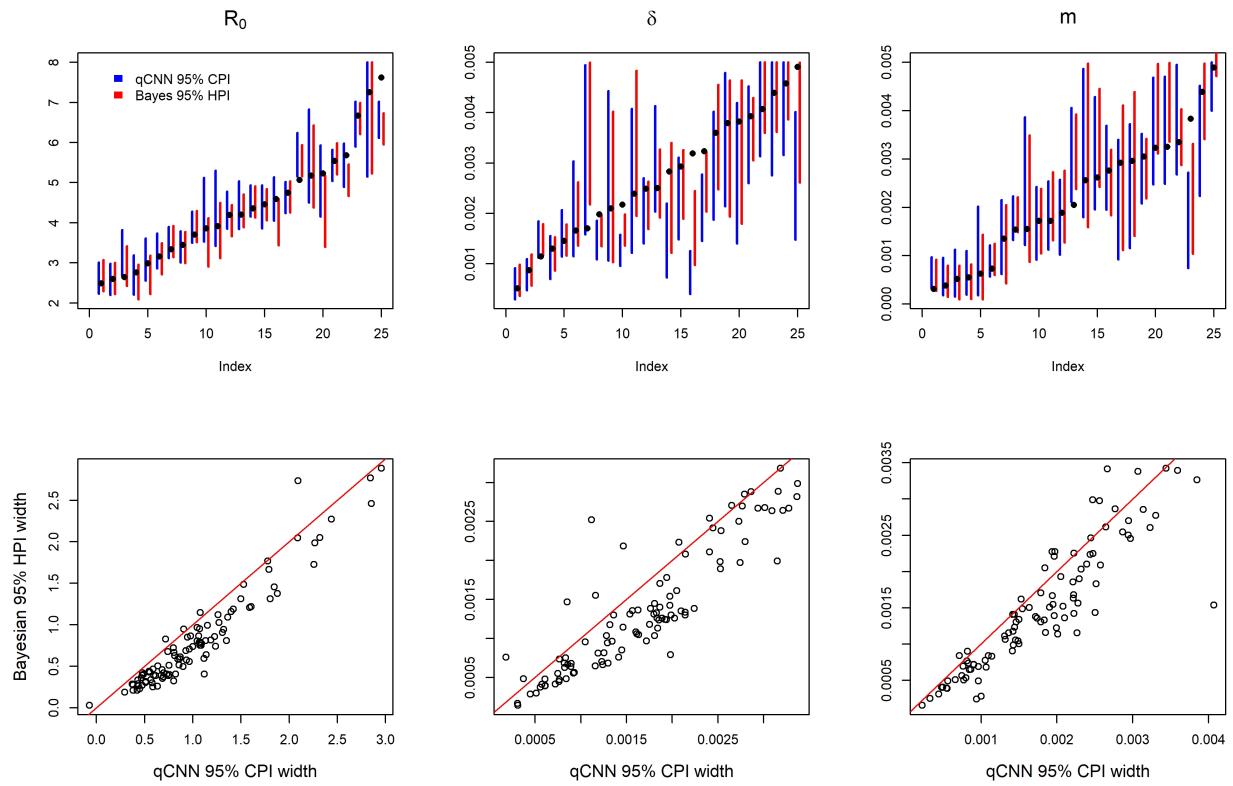


Figure S14: Comparison of CPI and HPI intervals for phylogeny error experiment. See SI Figure S4 for general details about plot.